# Supplementary Material for:

# Evaluating Large-Scale Propensity Score Performance

# Through Real-World and Synthetic Data Experiments

## Supplementary Material 1: Treatments - Anticoagulants Study

### Dabigatran New Users with Prior Atrial Fibrillation

Initial Event Cohort

People having any of the following:

- a drug era of dabigatran[4]

  - for the first time in the person's history

  - era start is on or after 2010-10-19

  - with age at era start $\geq 65$

with continuous observation of at least 183 days prior and 0 days after event index date, and limit initial events to: earliest event per person.

Inclusion Criteria 1: Has prior atrial fibrillation of atrial flutter diagnosis

Having any of the following criteria:

- at least 1 occurrences of a condition occurrence of Atrial fibrillation[2]
  starting between all days Before and 0 days After event index date

- or at least 1 occurrences of a condition occurrence of Atrial flutter[3]
  starting between all days Before and 0 days After event index date

Inclusion Criteria 2: Has no prior treatment with comparator drug (warfarin)

Having any of the following criteria:

- exactly 0 occurrences of a drug exposure of warfarin[13]
  starting between all days Before and 0 days Before event index date

1

Inclusion Criteria 3: Has no prior treatment with other anticoagulants (rivaroxaban or apixaban)

Having any of the following criteria:

- exactly 0 occurrences of a drug exposure of rivaroxaban[12]
  starting between all days Before and 0 days After event index date

- and exactly 0 occurrences of a drug exposure of apixaban[1]
  starting between all days Before and 0 days After event index date

Inclusion Criteria 4: Not in a skilled nursing facility or nursing home, or receiving hospice care on the index date

Having any of the following criteria:

- exactly 0 occurrences of a visit occurrence of long term care visit[10]
  starting between 0 days Before and 0 days After event index date

- and exactly 0 occurrences of a procedure of Hospice observations[9]
  starting between all days Before and 0 days After event index date

- and exactly 0 occurrences of an observation of Hospice observations[9]
  starting between all days Before and 0 days After event index date

Inclusion Criteria 5: Not undergoing dialysis or kidney transplant recipient

Having any of the following criteria:

- exactly 0 occurrences of a condition occurrence of Hemodialysis, peritoneal dialysis, or kidney transplant[7]
  starting between 183 days Before and 0 days After event index date

- and exactly 0 occurrences of a procedure of Hemodialysis, peritoneal dialysis, or kidney transplant[7]
  starting between 183 days Before and 0 days After event index date

- and exactly 0 occurrences of an observation of Hemodialysis, peritoneal dialysis, or kidney transplant[7]
  starting between 183 days Before and 0 days After event index date

Inclusion Criteria 6: No mitral valve disease, heart valve repair, or replacement in the prior 6 months

> Having any of the following criteria:

- exactly 0 occurrences of a condition occurrence of Heart valve disease, repair or replacement[6] starting between 183 days Before and 0 days After event index date

- and exactly 0 occurrences of a procedure of Heart valve disease, repair or replacement[6] starting between 183 days Before and 0 days After event index date

- and exactly 0 occurrences of an observation of Heart valve disease, repair or replacement[6] starting between 183 days Before and 0 days After event index date

Inclusion Criteria 7: No deep vein thrombosis or pulmonary embolism in the prior 6 months

> Having any of the following criteria:

- exactly 0 occurrences of a condition occurrence of Deep vein thrombosis[5] starting between 183 days Before and 0 days After event index date

- and exactly 0 occurrences of a condition occurrence of Pulmonary embolism[11] starting between 183 days Before and 0 days After event index date

Inclusion Criteria 8: No joint replacement surgery in the prior 6 months

> Having any of the following criteria:

- exactly 0 occurrences of a procedure of Hip/knee joint replacement or revision[8] starting between 183 days Before and 0 days After event index date

Cohort Exit Criteria

> Cohort end date is the end of the observation period that contains the index event.

**Warfarin New Users with Prior Atrial Fibrillation**

Initial Event Cohort

> People having any of the following:

- a drug era of warfarin[13]

  - for the first time in the person's history

  - era start is on or after 2010-10-19

  - with age at era start $\geq 65$

with continuous observation of at least 183 days prior and 0 days after event index date, and limit initial events to: earliest event per person.

Inclusion Criteria 1, 3-8: Same as Dabigatran Cohort

Inclusion Criteria 2: Has no prior treatment with comparator drug (dabigatran)

Having any of the following criteria:

- exactly 0 occurrences of a drug exposure of dabigatran[4]

  starting between all days Before and 0 days Before event index date

### Appendix: Concept Set Definitions

Codes given use the Observational Medical Outcomes Partnership Common Data Model Version 5 format (1). The "Vocabulary" column indicates the source vocabulary set for the concept. The "Excluded" column indicates whether the covariate is included (NO) or excluded (YES) from the cohort definition. The "Descendants" column incidates whether all descendent concepts in the vocabulary hierarchy are also incorporated.

1. apixaban

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 43013024 | apixaban | Drug | RxNorm | NO | YES |

2. Atrial fibrillation

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 313217 | Atrial fibrillation | Condition | SNOMED | NO | YES |

3. Atrial flutter

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 314665 | Atrial flutter | Condition | SNOMED | NO | YES |

## 4. dabigatran

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 40228152 | dabigatran etexilate | Drug | RxNorm | NO | YES |

## 5. Deep vein thrombosis

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 435887 | Antepartum deep vein thrombosis | Condition | SNOMED | YES | YES |
| 195562 | Hemorrhoids | Condition | SNOMED | YES | YES |
| 4179912 | Intracranial venous thrombosis | Condition | SNOMED | YES | YES |
| 318137 | Phlebitis and thrombophlebitis of intracranial sinuses | Condition | SNOMED | YES | YES |
| 199837 | Portal vein thrombosis | Condition | SNOMED | YES | YES |
| 438820 | Postpartum deep phlebothrombosis | Condition | SNOMED | YES | YES |
| 4235812 | Septic thrombophlebitis | Condition | SNOMED | YES | YES |
| 4187790 | Thrombosis of retinal vein | Condition | SNOMED | YES | YES |
| 318775 | Venous embolism | Condition | SNOMED | NO | YES |
| 444247 | Venous thrombosis | Condition | SNOMED | NO | YES |

## 6. Heart valve disease, repair or replacement

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 4060089 | H/O: artificial heart valve | Observation | SNOMED | NO | YES |
| 4195003 | Heart valve stenosis | Condition | SNOMED | NO | YES |
| 44782431 | History of mechanical heart valve replacement | Observation | SNOMED | NO | YES |
| 4013355 | Implantation of heart valve | Procedure | SNOMED | NO | YES |
| 4165384 | Implantation of heart valve prosthesis or synthetic device | Procedure | SNOMED | NO | YES |
| 2617335 | Md inr test revie inter mgmt | Observation | HCPCS | NO | YES |
| 43020459 | Mechanical breakdown of prosthetic heart valve | Condition | SNOMED | NO | YES |
| 312773 | Mechanical complication due to heart valve prosthesis | Condition | SNOMED | NO | YES |
| 4020159 | Mechanical complication of heart valve prosthesis | Condition | SNOMED | NO | YES |
| 44783274 | Mechanical heart valve replacement | Procedure | SNOMED | NO | YES |
| 315273 | Mitral valve stenosis | Condition | SNOMED | NO | YES |
| 4110937 | Non-rheumatic mitral valve stenosis | Condition | SNOMED | NO | YES |
| 2001447 | Open and other replacement of heart valve | Procedure | ICD9Proc | NO | YES |
| 2001448 | Open and other replacement of unspecified heart valve | Procedure | ICD9Proc | NO | YES |
| 4119522 | Prosthetic heart valve sample | Specimen | SNOMED | NO | YES |

| | | | | | |
|---|---|---|---|---|---|
| 4145884 | Prosthetic replacement of heart valve | Procedure | SNOMED | NO | YES |
| 2617334 | Provide inr test mater/equip | Observation | HCPCS | NO | YES |
| 4339971 | Reinsertion of heart valve, prosthetic | Procedure | SNOMED | NO | YES |
| 4121484 | Replacement of heart valve poppet, prosthetic | Procedure | SNOMED | NO | YES |
| 4013356 | Resuture of heart valve prosthesis, poppet | Procedure | SNOMED | NO | YES |
| 4181749 | Revision of prosthesis of heart valve | Procedure | SNOMED | NO | YES |
| 4304541 | Rheumatic mitral valve insufficiency AND aortic valve stenosis | Condition | SNOMED | NO | YES |

## 7. Hemodialysis, peritoneal dialysis, or kidney transplant

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 4126124 | Acute disorder of hemodialysis | Condition | SNOMED | NO | YES |
| 4092504 | Adequacy of hemodialysis | Observation | SNOMED | NO | YES |
| 435649 | Complication of hemodialysis | Condition | SNOMED | NO | YES |
| 40480136 | Dependence on hemodialysis | Observation | SNOMED | NO | YES |
| 4181476 | Dependence on hemodialysis due to end stage renal disease | Observation | SNOMED | NO | YES |
| 44786469 | Docrsn for cath maint dia | Observation | HCPCS | NO | YES |
| 4120120 | Hemodialysis | Procedure | SNOMED | NO | YES |
| 2101833 | Hemodialysis plan of care documented (esrd, p-esrd) | Observation | CPT4 | NO | YES |
| 4137616 | Hemodialysis-associated amyloidosis | Condition | SNOMED | NO | YES |
| 313232 | Hemodialysis-associated hypotension | Condition | SNOMED | NO | YES |
| 4300099 | Hemodialysis-associated pruritus | Condition | SNOMED | NO | YES |
| 4297919 | Hemodialysis-associated pseudoporphyria | Condition | SNOMED | NO | YES |
| 4297658 | Hemodialysis-associated secondary amyloidosis of skin | Condition | SNOMED | NO | YES |
| 4099603 | Megaloblastic anemia due to hemodialysis | Measurement | SNOMED | NO | YES |
| 44782924 | Misplacement of hemodialysis catheter | Condition | SNOMED | NO | YES |
| 44786470 | Patient receiving maintenance hemodialysis for greater than or equal to 90 days with a catheter as the mode of vascular access | Observation | HCPCS | NO | YES |
| 44786471 | Patient receiving maintenance hemodialysis for greater than or equal to 90 days without a catheter as the mode of vascular access | Observation | HCPCS | NO | YES |
| 43533281 | Patient receiving maintenance hemodialysis in an outpatient dialysis facility | Observation | HCPCS | NO | YES |
| 4324124 | Peritoneal dialysis | Procedure | SNOMED | NO | YES |
| 2003564 | Peritoneal dialysis | Procedure | ICD9Proc | NO | YES |
| 4300106 | Skin lesion associated with hemodialysis | Condition | SNOMED | NO | YES |
| 4046829 | Anesthesia for renal transplant, recipient | Procedure | SNOMED | NO | YES |

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 2109584 | Backbench reconstruction of cadaver or living donor renal allograft prior to transplantation; ureteral anastomosis, each | Procedure | CPT4 | NO | YES |
| 4021107 | Cadaveric renal transplant | Procedure | SNOMED | NO | YES |
| 4197300 | Donor renal transplantation | Procedure | SNOMED | NO | YES |
| 4324754 | Examination of recipient after kidney transplant | Procedure | SNOMED | NO | YES |
| 4002215 | Kidney implantation | Procedure | SNOMED | NO | YES |
| 4022805 | Live donor renal transplant | Procedure | SNOMED | NO | YES |
| 2003626 | Other kidney transplantation | Procedure | ICD9Proc | NO | YES |
| 40664909 | Patient receiving hemodialysis, peritoneal dialysis or kidney transplantation | Observation | HCPCS | NO | YES |
| 2109586 | Renal allotransplantation, implantation of graft; without recipient nephrectomy | Procedure | CPT4 | NO | YES |
| 2109589 | Renal autotransplantation, reimplantation of kidney | Procedure | CPT4 | NO | YES |
| 4163566 | Renal replacement | Procedure | SNOMED | NO | YES |
| 4346636 | Renal transplant arteriogram | Procedure | SNOMED | NO | YES |
| 4346505 | Renal transplant venogram | Procedure | SNOMED | NO | YES |
| 4347789 | Renal transplant venous sampling | Procedure | SNOMED | NO | YES |
| 2721092 | Simultaneous pancreas kidney transplantation | Procedure | HCPCS | NO | YES |
| 4322471 | Transplant of kidney | Procedure | SNOMED | NO | YES |
| 4343000 | Xenograft renal transplant | Procedure | SNOMED | NO | YES |

## 8. Hip/knee joint replacement or revision

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 2101660 | Anesthesia for open or surgical arthroscopic procedures on knee joint; total knee arthroplasty | Procedure | CPT4 | NO | YES |
| 2101635 | Anesthesia for open procedures involving hip joint; revision of total hip arthroplasty | Procedure | CPT4 | NO | YES |
| 2101634 | Anesthesia for open procedures involving hip joint; total hip arthroplasty | Procedure | CPT4 | NO | YES |
| 2104836 | Arthroplasty, acetabular and proximal femoral prosthetic replacement (total hip arthroplasty), with or without autograft or allograft | Procedure | CPT4 | NO | YES |
| 2103931 | Arthroplasty, elbow; with distal humerus and proximal ulnar prosthetic replacement (eg, total elbow) | Procedure | CPT4 | NO | YES |

| | | | | | |
|---|---|---|---|---|---|
| 2105103 | Arthroplasty, knee, condyle and plateau; medial AND lateral compartments with or without patella resurfacing (total knee arthroplasty) | Procedure | CPT4 | NO | YES |
| 2104837 | Conversion of previous hip surgery to total hip arthroplasty, with or without autograft or allograft | Procedure | CPT4 | NO | YES |
| 2104835 | Hemiarthroplasty, hip, partial (eg, femoral stem prosthesis, bipolar arthroplasty) | Procedure | CPT4 | NO | YES |
| 2000075 | Hip bearing surface, ceramic-on-ceramic | Procedure | ICD9Proc | NO | YES |
| 2000076 | Hip bearing surface, ceramic-on-polyethylene | Procedure | ICD9Proc | NO | YES |
| 2000074 | Hip bearing surface, metal-on-metal | Procedure | ICD9Proc | NO | YES |
| 2000073 | Hip bearing surface, metal-on-polyethylene | Procedure | ICD9Proc | NO | YES |
| 4001859 | Hip joint implantation | Procedure | SNOMED | NO | YES |
| 4134857 | Insertion of hip prosthesis | Procedure | SNOMED | NO | YES |
| 4207955 | Insertion of hip prosthesis, total | Procedure | SNOMED | NO | YES |
| 2005902 | Partial hip replacement | Procedure | ICD9Proc | NO | YES |
| 4162099 | Prosthetic arthroplasty of the hip | Procedure | SNOMED | NO | YES |
| 2000085 | Resurfacing hip, partial, acetabulum | Procedure | ICD9Proc | NO | YES |
| 2000084 | Resurfacing hip, partial, femoral head | Procedure | ICD9Proc | NO | YES |
| 2000083 | Resurfacing hip, total, acetabulum and femoral head | Procedure | ICD9Proc | NO | YES |
| 4010119 | Revision of hip replacement | Procedure | SNOMED | NO | YES |
| 2000070 | Revision of hip replacement, acetabular component | Procedure | ICD9Proc | NO | YES |
| 2000072 | Revision of hip replacement, acetabular liner and/or femoral head only | Procedure | ICD9Proc | NO | YES |
| 2000069 | Revision of hip replacement, both acetabular and femoral components | Procedure | ICD9Proc | NO | YES |
| 2000071 | Revision of hip replacement, femoral component | Procedure | ICD9Proc | NO | YES |
| 2000080 | Revision of knee replacement, femoral component | Procedure | ICD9Proc | NO | YES |
| 2000081 | Revision of knee replacement, patellar component | Procedure | ICD9Proc | NO | YES |
| 2000079 | Revision of knee replacement, tibial component | Procedure | ICD9Proc | NO | YES |
| 2000078 | Revision of knee replacement, total (all components) | Procedure | ICD9Proc | NO | YES |
| 45887894 | Revision of total hip arthroplasty | Procedure | CPT4 | NO | YES |

| 2104839 | Revision of total hip arthroplasty; acetabular component only, with or without autograft or allograft | Procedure | CPT4 | NO | YES |
|---|---|---|---|---|---|
| 2104838 | Revision of total hip arthroplasty; both components, with or without autograft or allograft | Procedure | CPT4 | NO | YES |
| 2104840 | Revision of total hip arthroplasty; femoral component only, with or without allograft | Procedure | CPT4 | NO | YES |
| 4266062 | Revision of total hip replacement | Procedure | SNOMED | NO | YES |
| 2105128 | Revision of femoral component of total arthroplasty of knee without allograft | Procedure | CPT4 | NO | YES |
| 2105129 | Revision of total knee arthroplasty, with or without allograft; femoral and entire tibial component | Procedure | CPT4 | NO | YES |
| 2000082 | Revision of total knee replacement, tibial insert (liner) | Procedure | ICD9Proc | NO | YES |
| 2005891 | Total hip replacement | Procedure | ICD9Proc | NO | YES |
| 2005904 | Total knee replacement | Procedure | ICD9Proc | NO | YES |
| 4203771 | Total replacement of hip | Procedure | SNOMED | NO | YES |
| 2005903 | Revision of hip replacement, not otherwise specified | Procedure | ICD9Proc | NO | YES |
| 2104914 | Open treatment of femoral fracture, proximal end, neck, internal fixation or prosthetic replacement | Procedure | CPT4 | YES | YES |

## 9. Hospice Observations

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 40483762 | Acute care hospice service | Observation | SNOMED | NO | YES |
| 4123927 | Admission to hospice | Observation | SNOMED | NO | YES |
| 4086294 | Admission to hospice for respite | Observation | SNOMED | NO | YES |
| 4137269 | Discharge from hospice | Observation | SNOMED | NO | YES |
| 4137272 | Discharge from hospice day hospital | Observation | SNOMED | NO | YES |
| 4062333 | Full care by hospice | Observation | SNOMED | NO | YES |
| 40481548 | Home hospice service | Observation | SNOMED | NO | YES |
| 4109386 | Hospice | Observation | SNOMED | NO | YES |
| 4301458 | Hospice care | Observation | SNOMED | NO | YES |
| 4301459 | Hospice care management | Procedure | SNOMED | NO | YES |
| 2720815 | Hospice care provided in inpatient hospice facility | Observation | HCPCS | NO | YES |
| 2720814 | Hospice care provided in inpatient hospital | Observation | HCPCS | NO | YES |
| 2720817 | Hospice care provided in inpatient psychiatric facility | Observation | HCPCS | NO | YES |
| 2720816 | Hospice care provided in long term care facility | Observation | HCPCS | NO | YES |

| | | | | | |
|---|---|---|---|---|---|
| 2720812 | Hospice care provided in nursing long term care facility (ltc) or non-skilled nursing facility (nf) | Observation | HCPCS | NO | YES |
| 2720813 | Hospice care provided in skilled nursing facility (snf) | Observation | HCPCS | NO | YES |
| 2617270 | Hospice care supervision | Observation | HCPCS | NO | YES |
| 2721445 | Hospice care, in the home, per diem | Observation | HCPCS | NO | YES |
| 2721700 | Hospice continuous home care; per hour | Observation | HCPCS | NO | YES |
| 2721702 | Hospice general inpatient care; per diem | Observation | HCPCS | NO | YES |
| 40664432 | Hospice home care provided in a hospice facility | Observation | HCPCS | NO | YES |
| 2721701 | Hospice inpatient respite care; per diem | Observation | HCPCS | NO | YES |
| 2721703 | Hospice long term care, room and board only; per diem | Observation | HCPCS | NO | YES |
| 2720811 | Hospice or home health care provided in assisted living facility | Observation | HCPCS | NO | YES |
| 38003372 | Hospice Room Board-Nursing facility | Revenue Code | Revenue Code | NO | YES |
| 2721699 | Hospice routine home care; per diem | Observation | HCPCS | NO | YES |
| 38003368 | Hospice Service - Continuous Home Care | Revenue Code | Revenue Code | NO | YES |
| 38003366 | Hospice Service - General Classification | Revenue Code | Revenue Code | NO | YES |
| 38003370 | Hospice Service - General Inpatient Care (Non-respite) | Revenue Code | Revenue Code | NO | YES |
| 38003369 | Hospice Service - Impatient Respite Care | Revenue Code | Revenue Code | NO | YES |
| 38003373 | Hospice Service - Other | Revenue Code | Revenue Code | NO | YES |
| 38003371 | Hospice Service - Physician Services | Revenue Code | Revenue Code | NO | YES |
| 38003367 | Hospice Service - Routine Home Care | Revenue Code | Revenue Code | NO | YES |
| 38003131 | Incremental Nursing Charge Rate - Hospice | Revenue Code | Revenue Code | NO | YES |
| 38003066 | Private (Deluxe) - Hospice | Revenue Code | Revenue Code | NO | YES |
| 38003036 | Room Board - Private (Medical or General) - Hospice | Revenue Code | Revenue Code | NO | YES |
| 38003046 | Room Board - Semi-private Two Bed (Medical or General) - Hospice | Revenue Code | Revenue Code | NO | YES |

| 38003076 | Room Board Ward (Medical or General) - Hospice | Revenue Code | Revenue Code | NO | YES |
|---|---|---|---|---|---|
| 4082084 | Routine admission to hospice | Observation | SNOMED | NO | YES |
| 4140947 | Seen in hospice | Observation | SNOMED | NO | YES |
| 38003056 | Semi-Private - Three and Four Beds - Hospice | Revenue Code | Revenue Code | NO | YES |
| 915618 | Services performed by care coordinator in the hospice setting, each 15 minutes | Observation | HCPCS | NO | YES |
| 915614 | Services performed by chaplain in the hospice setting, each 15 minutes | Observation | HCPCS | NO | YES |
| 915615 | Services performed by dietary counselor in the hospice setting, each 15 minutes | Observation | HCPCS | NO | YES |
| 915616 | Services performed by other counselor in the hospice setting, each 15 minutes | Observation | HCPCS | NO | YES |
| 915619 | Services performed by other qualified therapist in the hospice setting, each 15 minutes | Observation | HCPCS | NO | YES |
| 915620 | Services performed by qualified pharmacist in the hospice setting, each 15 minutes | Observation | HCPCS | NO | YES |
| 915617 | Services performed by volunteer in the hospice setting, each 15 minutes | Observation | HCPCS | NO | YES |
| 4062044 | Shared care - hospice and GP | Observation | SNOMED | NO | YES |
| 2514512 | Supervision of care of hospice patient, without patient present - 15-29 minutes | Procedure | CPT4 | NO | YES |
| 4086777 | Urgent admission to hospice | Observation | SNOMED | NO | YES |

## 10. long term care visit

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 42898160 | Long Term Care Visit | Visit | Visit | NO | YES |

## 11. Pulmonary embolism

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 40480461 | Infarction of lung due to iatrogenic pulmonary embolism | Condition | SNOMED | NO | YES |
| 435026 | Obstetric pulmonary embolism | Condition | SNOMED | YES | YES |
| 440417 | Pulmonary embolism | Condition | SNOMED | NO | YES |
| 40479606 | Septic pulmonary embolism | Condition | SNOMED | NO | YES |

## 12. rivaroxaban

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 40241331 | rivaroxaban | Drug | RxNorm | NO | YES |

13. warfarin

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 1310149 | Warfarin | Drug | RxNorm | NO | YES |

# Supplementary Material 2: Outcome - Anticoagulants Study

## Incident Intracranial Hemorrhage, Observed in Inpatient Setting

## Initial Event Cohort

People having any of the following:

- a condition occurrence of Intracranial hemorrhage[1]

    ○ condition type is any of: Inpatient detail - primary, Inpatient header - primary, Primary Condition, Inpatient detail - 1st position, Inpatient header - 1st position

    ○ visit occurrence is any of: Emergency Room Visit, Inpatient Visit

with continuous observation of at least 0 days prior and 0 days after event index date, and limit initial events to: earliest event per person.

1. Intracranial hemorrhage

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 376713 | Cerebral hemorrhage | Condition | SNOMED | NO | YES |
| 4014781 | Closed traumatic subdural hemorrhage | Condition | SNOMED | NO | YES |
| 252477 | Extradural hemorrhage following injury without open intracranial wound | Condition | SNOMED | NO | YES |
| 439847 | Intracranial hemorrhage | Condition | SNOMED | NO | NO |
| 42873157 | Intracranial hemorrhage following injury | Condition | SNOMED | NO | NO |
| 436430 | Nontraumatic extradural hemorrhage | Condition | SNOMED | NO | YES |
| 432923 | Subarachnoid hemorrhage | Condition | SNOMED | NO | YES |
| 435959 | Subarachnoid hemorrhage following injury without open intracranial wound | Condition | SNOMED | NO | YES |
| 439040 | Subdural hemorrhage | Condition | SNOMED | NO | YES |
| 260841 | Perinatal subarachnoid hemorrhage | Condition | SNOMED | YES | YES |
| 436519 | Perinatal intraventricular hemorrhage | Condition | SNOMED | YES | YES |
| 4071732 | Intracranial nontraumatic hemorrhage of fetus and newborn | Condition | SNOMED | YES | YES |
| 4345688 | Intracerebral hemorrhage in fetus or newborn | Condition | SNOMED | YES | YES |

| 443752 | Ventricular hemorrhage | Condition | SNOMED | YES | YES |
|---|---|---|---|---|---|
| 4174299 | Perinatal intracranial hemorrhage | Condition | SNOMED | YES | YES |
| 380113 | Hemorrhage in optic nerve sheaths | Condition | SNOMED | YES | YES |
| 42872434 | Intracranial hematoma | Condition | SNOMED | NO | YES |

# Supplementary Material 3: Treatments - NSAIDs Study

## Celecoxib New Users

Initial Event Cohort

People having any of the following:

- a drug era of celecoxib[2]

  ○ for the first time in the person's history

  ○ era start is between 2004-01-01 and 2007-12-31 (inclusive)

  ○ with age at era start $\geq 16$

with continuous observation of at least 365 days prior and 0 days after event index date, and limit initial events and qualifying cohort to: earliest event per person.

Inclusion Criteria:

Having all of the following criteria:

- exactly 0 occurrences of a condition occurrence of Broad malignancies excluding skin cancer (including primary, secondary)[1]

  starting between 365 days Before and 0 days Before event index date

- and exactly 0 occurrences of a drug exposure of NSAIDs[4]

  starting between 365 days Before and 1 days Befor event index date

Cohort Exit Criteria

Cohort definition end date is index event's end date plus 0 days

**Diclofenac New Users**

Initial Event Cohort

People having any of the following:

- a drug era of diclofenac[3]

    ○ for the first time in the person's history

    ○ era start is between 2004-01-01 and 2007-12-31 (inclusive)

    ○ with age at era start $\geq 16$

with continuous observation of at least 183 days prior and 0 days after event index date, and limit initial events to: earliest event per person.

Inclusion Criteria: Same as celecoxib cohort

Cohort Exit Criteria: Same as celecoxib cohort

**Appendix: Concept Set Definitions**

1. Broad malignancies excluding skin cancer (including primary, secondary)

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 4179980 | Malignant basal cell neoplasm of skin | Condition | SNOMED | YES | YES |
| 443392 | Malignant neoplastic disease | Condition | SNOMED | NO | YES |
| 4300118 | Squamous cell carcinoma | Condition | SNOMED | YES | YES |

2. celecoxib

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 1118084 | celecoxib | Drug | RxNorm | NO | YES |

3. diclofenac

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|
| 1124300 | Diclofenac | Drug | RxNorm | NO | YES |

4. NSAIDs

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
|---|---|---|---|---|---|

| 21603933 | ANTIINFLAMMATORY AND AN-TIRHEUMATIC PRODUCTS, NON-STEROIDS | Drug | ATC | NO | YES |
| --- | --- | --- | --- | --- | --- |

# Supplementary Material 4: Outcome - NSAIDs Study

## Upper gastrointestinal complication (UGIC) events

Initial Event Cohort

People having any of the following:

- a condition era of Upper gastrointestinal complication (UGIC)[1]

with continuous observation of at least 0 days prior and 0 days after event index date, and limit initial events and qualifying cohort to: earliest event per person.

Inclusion Criteria:

Having all of the following criteria:

- exactly 0 occurrences of a condition occurrence of Upper gastrointestinal complication (UGIC)[1] starting between 30 days Before and 1 day Before event index date

### Appendix: Concept Set Definitions

1. Broad malignancies excluding skin cancer (including primary, secondary)

| Concept Id | Concept Name | Domain | Vocabulary | Excluded | Descendants |
| --- | --- | --- | --- | --- | --- |
| 192671 | Gastrointestinal hemorrhage | Condition | SNOMED | NO | YES |
| 194158 | Perinatal gastrointestinal hemorrhage | Condition | SNOMED | YES | YES |
| 26441 | Bleeding ulcer of esophagus | Condition | SNOMED | YES | YES |
| 316457 | Mallory-Weiss syndrome | Condition | SNOMED | YES | YES |
| 46273478 | Rectal hemorrhage due to ulcerative colitis | Condition | SNOMED | YES | YES |
| 4071070 | Neonatal hematemesis | Condition | SNOMED | YES | YES |
| 4048286 | Neonatal rectal hemorrhage | Condition | SNOMED | YES | YES |
| 4103703 | Melena | Condition | SNOMED | NO | YES |
| 4048608 | Neonatal melena | Condition | SNOMED | YES | YES |
| 4265600 | Gastric ulcer | Condition | SNOMED | NO | YES |
| 4027663 | Peptic ulcer | Condition | SNOMED | NO | YES |
| 4059178 | Gastrojejunal ulcer | Condition | SNOMED | NO | YES |
| 4198381 | Duodenal ulcer disease | Condition | SNOMED | NO | YES |
| 4172869 | Peptic ulcer of newborn | Condition | SNOMED | YES | YES |

# Supplementary Material 5: Synthetic Framework

## Notation

There are $N$ total study subjects indexed by $i$, and each subject has have treatment indicator $w_i$ and $p$-length baseline covariate vector $\boldsymbol{x}_i$. The event time is $t_i$, and $\delta_i$ is the censoring indicator, with $\delta_i = 0$ indicating censoring and $\delta_i = 1$ indicating the outcome of interest. Under the proportional hazards model, $\eta$ and $\boldsymbol{\beta}$ are the log hazard ratios for the treatment and the baseline covariates, respectively; the subject-specific hazard is then $\theta_i = w_i\eta + \boldsymbol{x}_i\boldsymbol{\beta}$. The baseline survival function $S(t)$ traces the probability of surviving to time $t$ after treatment initiation and The analagous baseline censoring function is $C(t)$.

## Estimate simulation components

Outcome simulation requires estimates for $S(t)$, $C(t)$, and $\boldsymbol{\beta}$. We estimate $S(t)$ by fitting a distribution to the observed outcome of interest, and $C(t)$ by fitting a distribution to the censoring times. Critically, the censoring function must be covariate-free to maintain non-informative censoring for the proportional hazards model, meaning that a subject's censoring time and survival time are independent. This point is overlooked in the "plasmode" framework, leading to inaccurate true hazard ratios that are not proportional hazards. Possible forms for $S(t)$ and $C(t)$ include parametric distributions such as exponential, Weibull, Gompertz, gamma, and lognormal; discrete nonparametric estimators such as the Breslow and Kalbfleisch-Prentice estimators (2) (which without covariates are respectively the Nelson-Aalen and Kaplan-Meier estimators); and nonparametric spline functions (3). The $S(t)$ distribution determines how the covariate coefficients $\boldsymbol{\beta}$ are estimated. For parametric and spline estimators, the parameters that characterize $S(t)$ are jointly estimated with the covariate coefficients often using maximum likelihood estimation on the full survival likelihood function. For the discrete nonparametric estimators, covariate coefficients are first estimated via the partial likelihood function, and then used to produce $S(t)$ (4, 5). The subject-specific hazard is then $\theta_i = w_i\hat{\eta} + \boldsymbol{x_i}\hat{\boldsymbol{\beta}}$ and the subject-specific survival function $S(t)^{\exp\{\theta_i\}}$, where $\hat{\eta}$ and $\hat{\boldsymbol{\beta}}$ are maximum likelihood estimates.

The nonparametric Breslow and Kalbfleisch-Prentice baseline survival function estimators are discrete functions that assign outcome probability only to empirically observed outcome times. In

scenarios with relatively extreme hazard ratios or outcome prevalences, these discrete estimators can affect estimation bias by simulating excessive outcome time ties. To avoid this, we smooth these estimators to produce a non-disjoint function by linearly connecting each point in the estimators. Let the $M$ increasingly ordered observed outcome times be $t_{(k)}, 1 \leq k \leq M$. Assume $S(t_{(0)}) = 1$ and $S(t_{(M+1)}) = 0$. We assign to $S(t : t_{(k)} < t < t_{(k+1)})$ the value at $t$ on the straight line between $(t_{(k)}, S(t_{(k)}))$ and $(t_{(k+1)}, S(t_{(k+1)}))$. Because only the ordering of survival times affects the proportional hazards likelihood, our simple smoothing approach suffices to prevent excessively tied times.

For our simulations, we obtain $\hat{\boldsymbol{\beta}}$ through partial likelihood maximum likelihood estimation, and include $L_1$-regularization on all covariates except treatment to promote model fitting (6). We manually select the regularization penalty to yield an approximate model size of 500, coinciding with the number of covariates selected by the hdPS. We use the Breslow estimator for $S(t)$ (5), and the Nelson-Aalen estimator for $C(t)$.

**Simulate outcome and censoring times**

Under the proportional hazards framework, each subject's survival process is $S(t)^{\exp\{\theta_i\}}$ and censoring process $C(t)$. We use inverse transform sampling and draw for each subject two $\text{Unif}(0,1)$ random variables $R_{i,s}$ and $R_{i,c}$. The respective outcome and censoring times are $t_{i,s} = \min\{t : S_i(t) \leq R_{i,s}\}$ and $t_{i,c} = \min\{t : C(t) \leq R_{i,c}\}$. The final simulated event is the miminum time:

$$t_i = \min\{t_{i,s}, t_{i,c}\}$$

$$\delta_i = \begin{cases} 1 & t_{i,s} < t_{i,c} \\ 0 & t_{i,s} \geq t_{i,c} \end{cases}. \tag{1}$$

**Adjust simulation for hazard ratio and outcome prevalence**

To simulate under a desired treatment hazard ratio $\eta^*$, we replace the empirically estimated $\hat{\eta}$ by $\eta^*$: $\theta_i = w_i \eta^* + x_i \hat{\beta}$. The expected resultant simulated outcome prevalence (OP) is

$$p = \frac{1}{N} \sum_i \int_0^\infty \Pr\left(t_{i,s} = t < t_{i,c}\right) \, \mathrm{d}t = \frac{1}{N} \sum_i \int_0^\infty \left(\frac{\partial}{\partial t} S(t)^{\exp\{\theta_i\}}\right) C(t) \, \mathrm{d}t. \tag{2}$$

Let $t_{(k)}$ be the observed outcome times; the corresponding equation for discrete estimators is

$$p = \frac{1}{N} \sum_i \sum_{t_{(k)}} \left[ S\left(t_{(k-1)}\right)^{\exp\{\theta_i\}} - S\left(t_{(k)}\right)^{\exp\{\theta_i\}} \right] C(t_{(k)}). \tag{3}$$

To adjust the simulated outcome prevalence, several different approaches are available that can suit the investigator. Each requires solving outcome prevalence Equation (2) for an adjustment factor $\gamma$:

- Adjust outcome hazard – multiply the outcome baseline hazard by a fixed constant $\gamma \in (0, \infty)$ so the baseline survival function becomes $S(t) \to S(t)^\gamma$. This aligns with a proportional hazards interpretation of adjusting the baseline function.

- Adjust censoring hazard – multiply the censoring hazard by a fixed constant $\gamma \in (0, \infty)$, so the censoring function becomes $C(t) \to C(t)^\gamma$. Note that the maximum achievable outcome prevalence is $1 - S(T)$ when $C(t) = 1$, where $T$ is the maximum observed outcome time.

- Accelerate survival process – scale time by a factor $\gamma \in (0, \infty)$ in the survival function, to obtain a new survival function $S(t) \to S(\gamma t)$. This aligns with an accelerated life interpretation of adjusting the baseline function

For our simulations, we use the first approach, similar to the "plasmode" framework.

## Supplementary Material 6: Anticoagulants Study Negative Controls

Negative controls were selected using the following criteria from (7):

- No evidence found in literature on clinical trials using the method proposed by Avillach (8)

- No evidence found in literature using the method used in SemMedDB (9)

- No evidence found in the structured product label (US and EU).

- FAERS Proportional Reporting Ratio (PRR) needed to be less than 2.

Negative controls were rank-ordered by prevalence in study cohort, and manually reviewed until 50 controls were selected. Negative controls with fewer than 0.02% prevalence were discarded.

| | |
|---|---|
| Acute bronchitis | Allergic rhinitis |
| Anxiety disorder | Arthritis of spine |
| Arthropathy of knee joint | Atelectasis |
| Barrett's esophagus | Blepharitis |
| Bronchiectasis | Bundle branch block |
| Cellulitis | Chronic sinusitis |
| Chronic ulcer of skin | Communication disorder |
| Crohn's disease | Curvature of spine |
| Cutis laxa | Diabetic renal disease |
| Diabetic retinopathy | Dislocation of joint |
| Dyssomnia | Dysuria |
| Effusion of joint | Fracture of upper limb |
| Gallstone | Gammopathy |
| Human papilloma virus infection | Hyperplasia of prostate |
| Inflammation of sacroiliac joint | Ingrowing nail |
| Malignant tumor of breast | Multiple sclerosis |
| Neck pain | Neurologic disorder associated with diabetes mellitus |
| Obesity | Osteomyelitis |
| Otitis media | Peripheral vertigo |
| Plantar fasciitis | Presbyopia |
| Prolapse of female genital organs | Psychotic disorder |
| Seborrheic keratosis | Simple goiter |
| Sleep apnea | Superficial mycosis |
| Urge incontinence of urine | Urinary tract infectious disease |
| Verruca vulgaris | |

# Supplementary Material 7: NSAIDs Study Negative controls

Negative controls selected similarly to Anticoagulants study (Supplementary Material 6).

| | |
|---|---|
| Adjustment disorder | Alcohol abuse |
| Aphasia | Astigmatism |
| Atelectasis | Bell's palsy |
| Candida infection of genital region | Carcinoma in situ of breast |
| Chalazion | Curvature of spine |
| Cutis laxa | Deficiency of macronutrients |
| Diabetic oculopathy | Drug withdrawal |
| Exostosis | Fibrocystic disease of breast |
| Human papilloma virus infection | Hydrocele |
| Ingrowing nail | Intracranial injury |
| Meniere's disease | Non-toxic nodular goiter |
| Non-toxic uninodular goiter | Presbyopia |
| Scoliosis deformity of spine | Secondary malignant neoplastic disease |
| Tinea pedis | Verruca vulgaris |
| Vitamin D deficiency | |

# Supplementary Material 8: Covariate Sets

## OHDSI Covariates

- Demographics (age in 5-year increments, gender, year of index date)

- Condition Occurrence (condition occurrence in lookback window)

    - in 365 days prior to index date

    - in 180 days prior to index date

    - in 30 days prior to index date

- Condition Era (span of time when person assumed to have condition)

- ever

- overlapping with index date

- Drug Exposure (drug occurrence in lookback window)

  - in 365 days prior to index date

  - in 30 days prior to index date

- Drug Era (span of time when person assumed to have drug)

  - in 365 days prior to index date

  - in 30 days prior to index date

  - ever

  - overlapping with index date

- Procedure Occurrence

  - in 365 days prior to index date

  - in 30 days prior to index date

- Observation

  - in 365 days prior to index date

  - in 30 days prior to index date

  - observations count in 365 days prior to index date

- Measurement

  - in 365 days prior to index date

  - in 30 days prior to index date

  - high measurement in 180 days prior to index date

  - low measurement in 180 days prior to index date

  - measurements count in 365 days prior to index date

- Concepts Count in 365 days prior to index date

- Risk Scores (Charlson, DCSI, CHADS2)

**hdPS Covariates**

- Demographics (age in 5-year increments, gender, year of index date)

- Inpatient 3-Digit ICD9 Condition Codes

  ○ in 180 days prior to index date

  ○ $> 50^{\text{th}}$ percentile in 180 days prior to index date

  ○ $> 75^{\text{th}}$ percentile in 180 days prior to index date

- Outpatient 3-Digit ICD9 Condition Codes

  ○ in 180 days prior to index date

  ○ $> 50^{\text{th}}$ percentile in 180 days prior to index date

  ○ $> 75^{\text{th}}$ percentile in 180 days prior to index date

- Drug Ingredient Exposure

  ○ in 180 days prior to index date

  ○ $> 50^{\text{th}}$ percentile in 180 days prior to index date

  ○ $> 75^{\text{th}}$ percentile in 180 days prior to index date

- Inpatient CPT-4 Procedure Codes

  ○ in 180 days prior to index date

  ○ $> 50^{\text{th}}$ percentile in 180 days prior to index date

  ○ $> 75^{\text{th}}$ percentile in 180 days prior to index date

- Outpatient CPT-4 Proceudre Codes

  ○ in 180 days prior to index date

  ○ $> 50^{\text{th}}$ percentile in 180 days prior to index date

  ○ $> 75^{\text{th}}$ percentile in 180 days prior to index date

**Covariates Excluded in Anticoagulants Study**

RxNorm drugs Dabigatran (concept id 40228152), Warfarin (concept id 1310149); and all decescendent drugs / formulations.

**Covariates Excluded in NSAIDs Study**

RxNorm drugs Diclofenac (concept id 1124300), Celecoxib (concept id 1118084), Ingredients NSAIDs (concept id 21603933), Coxibs (concept id 21603991); and all decescendent drugs / formulations.

## Supplementary Material 9: Propensity Score Methods

We use the default hdPS settings including considering only the 200 most prevalent covariates in each "data dimension" and selecting the top 500 overall ranked covariates.

After estimating the PS, we perform PS matching and then estimate the treatment hazard ratio using a stratified Cox survival outcome model with treatment as the only covariate. We avoid one-to-one matching due to inferior covariate balance (10) and bias reduction (11), and instead use variable length matching (12) with a maximum ratio of 10:1 and a propensity score caliper of 0.05, and use a greedy matching algorithm (13). We use the less prevalent treatment as the "one" in the many-to-one matching to maximize the number of subjects that are matched.

## Supplementary Material 10: Propensity Score Estimate Existence

Supplementary Table 1 shows the convergence frequency of the hdPS across a spectrum of simulated sample sizes and outcome prevalences. The NSAIDs study demonstrates poorer hdPS estimate existence compared to the Anticoagulants study, including a complete failure of the exposure-based hdPS at any cohort size. In both studies, simulations with smaller cohorts and lower outcome prevalences have less likely PS estimate existence. Overall, the feasible use of the hdPS without regularization is data dependent and better suited for larger cohorts and more common outcomes. The inclusion of statistical regularization in the PS model readily solves this convergence problem.

| Study | Sample size | exp-hdPS | bias-hdPS | | | |
|---|---|---|---|---|---|---|
| | | | 0.5% | 1% | 5% | 10% |
| Anticoagulants | 72,489 | 100%* | 100% | 100% | 100% | 100% |
| | 30,000 | 100% | 99% | 100% | 100% | 100% |
| | 10,000 | 76% | 63% | 83% | 100% | 99% |
| | 5,000 | 26% | 46% | 24% | 80% | 81% |
| NSAIDs | 121,317 | 0%* | 99% | 99% | 100% | 100% |
| | 30,000 | 0% | 0% | 9% | 43% | 63% |
| | 10,000 | 0% | 1% | 0% | 0% | 0% |
| | 5,000 | 0% | 1% | 0% | 0% | 0% |

Supplementary Table 1: Frequency of propensity score estimate existence "convergence" in simulations using hdPS without regularization, over 100 simulations. Smaller cohort sizes are drawn repeatedly without replacement from the full cohort size. *exp-hdPS on full cohorts is tested only once.

## Supplementary Material 11: Bias Towards Null in Simulation Experiment

In the simulation experiments, the data are simulated under a proportional hazards survival model, matched on estimated propensity scores, and then fit under a Cox proportional hazards outcome model with treatment as the only covariate. The partial likelihood of the outcome model is

$$\log L = \sum_k \sum_{i:\delta_{i,k}=1} \left( w_{i,k}\eta - \log \sum_{j \in R_{i,k}} e^{w_{j,k}\eta} \right), \tag{4}$$

where $k$ indexes the strata/matched sets, $i$ indexes the subjects in each strata, $w_{i,k}$ is the exposure indicator, $\delta_{i,k}$ is the censoring indicator, and $R_{i,k}$ is the risk set for a subject in the strata.

We consider the case of one-to-one matching, in which the maximum likelihood estimate can be derived analytically. In matched set $k$, let $t_{0,k}$ and $t_{1,k}$ be the outcome times of the untreated and treated subjects, respectively. Because there is contribution to the log-likelihood only from non-censored subjects, the only configurations for matched sets that contribute to the log-likelihood are as follows:

| Case | $\delta_{0,k}$ | $\delta_{1,k}$ | $t_{0,k}$ ? $t_{1,k}$ | log-L contribution |
|---|---|---|---|---|
| A | 1 | 0 | $\leq$ | $0 - \log(1 + e^{\eta})$ |
| B | 0 | 1 | $\geq$ | $\eta - \log(1 + e^{\eta})$ |
| C | 1 | 1 | $<$ | $0 - \log(1 + e^{\eta})$ |
| D | 1 | 1 | $>$ | $\eta - \log(1 + e^{\eta})$ |
| E | 1 | 1 | $=$ | $\eta - 2\log(1 + e^{\eta})$ |

Supplementary Table 2: Contribution of matched sets to likelihood

Let $N_A$ denote the number of sets with case A, and likewise for the other cases. The log-likelihood can be written:

$$\log L = (N_A + N_C + N_E)[-\log(1 + e^\eta)] + (N_B + N_D + N_E)[\eta - \log(1 + e^\eta)]. \tag{5}$$

The sum $N_A + N_C + N_E$ is equal to the number of sets with $\delta_{0,k} = 1$ and $t_{0,k} \leq t_{1,k}$. Similarly, the sum $N_B + N_D + N_E$ is equal to the number of sets with $\delta_{1,k} = 1$ and $t_{0,k} \geq t_{1,k}$. Define $N_0 = N_A + N_C + N_E$ and $N_1 = N_B + N_D + N_E$. A second-derivative test with respect to $\eta$ shows this is concave in $\eta$; we solve by setting the first-derivative to zero:

$$
\begin{aligned}
\frac{\partial \log L}{\partial \eta}\bigg|_{\eta = \hat\eta} = N_0\left(-\frac{e^{\hat\eta}}{1 + e^{\hat\eta}}\right) + N_1\left(1 - \frac{e^{\hat\eta}}{1 + e^{\hat\eta}}\right) = 0 \\
N_0(-e^{\hat\eta}) + N_1(1) = 0 \\
e^{\hat\eta} = \frac{N_1}{N_0} \\
\hat\eta = \log N_1 - \log N_0.
\end{aligned} \tag{6}
$$

The observed bias to the null arises from the nonlinear effect that differences in hazard under the generative outcome model between matched subjects has on the estimated hazard ratio. In our simulation process, each subject has a simulated outcome and censoring time, and the minimum is the observed time. Let $t_{i,k,s}$ and $t_{i,k,c}$ be the simulated outcome and censoring time for subject $i$ in set $k$. Then $N_0$ is a sum of indicators $N_0 = \sum_k \mathbb{1}_{\{t_{0,k,s} \leq \min\{t_{0,k,c}, t_{1,k,s}, t_{1,k,c}\}\}}$ and $N_1$ is a sum of indicators $N_1 = \sum_k \mathbb{1}_{\{t_{1,k,s} \leq \min\{t_{1,k,c}, t_{0,k,s}, t_{0,k,c}\}\}}$. These indicators have the probabilities:

$$
\begin{aligned}
\Pr\left(t_{0,k,s} \leq \min\{t_{0,k,c}, t_{1,k,s}, t_{1,k,c}\}\right) = \int_0^\infty \left(\frac{\partial}{\partial t} S(t)^{\exp\{\theta_{0,k}\}}\right) S(t)^{\exp\{\theta_{1,k}\}} C(t)C(t) \mathrm{d}t \text{ and} \\
\Pr\left(t_{1,k,s} \leq \min\{t_{1,k,c}, t_{0,k,s}, t_{0,k,c}\}\right) = \int_0^\infty \left(\frac{\partial}{\partial t} S(t)^{\exp\{\theta_{1,k}\}}\right) S(t)^{\exp\{\theta_{0,k}\}} C(t)C(t) \mathrm{d}t.
\end{aligned} \tag{7}
$$

The subject-specific hazards are $\theta_{0,k} = \boldsymbol{x_{0,k}\beta}$ and $\theta_{1,k} = \eta^* + \boldsymbol{x_{1,k}\beta}$. There is a complicated relationship between $\eta^*$ and $\hat\eta$ that contributes to the observed bias.

We use simple simulations to empirically reproduce the bias towards the null observed in our simulation experiment results. We assign hazards to subjects in 1-1 matched sets, simulate event times, and then fit a stratified Cox model. We draw $\theta_{0,k}$ and $\theta_{1,k}$ from independent normal dis-
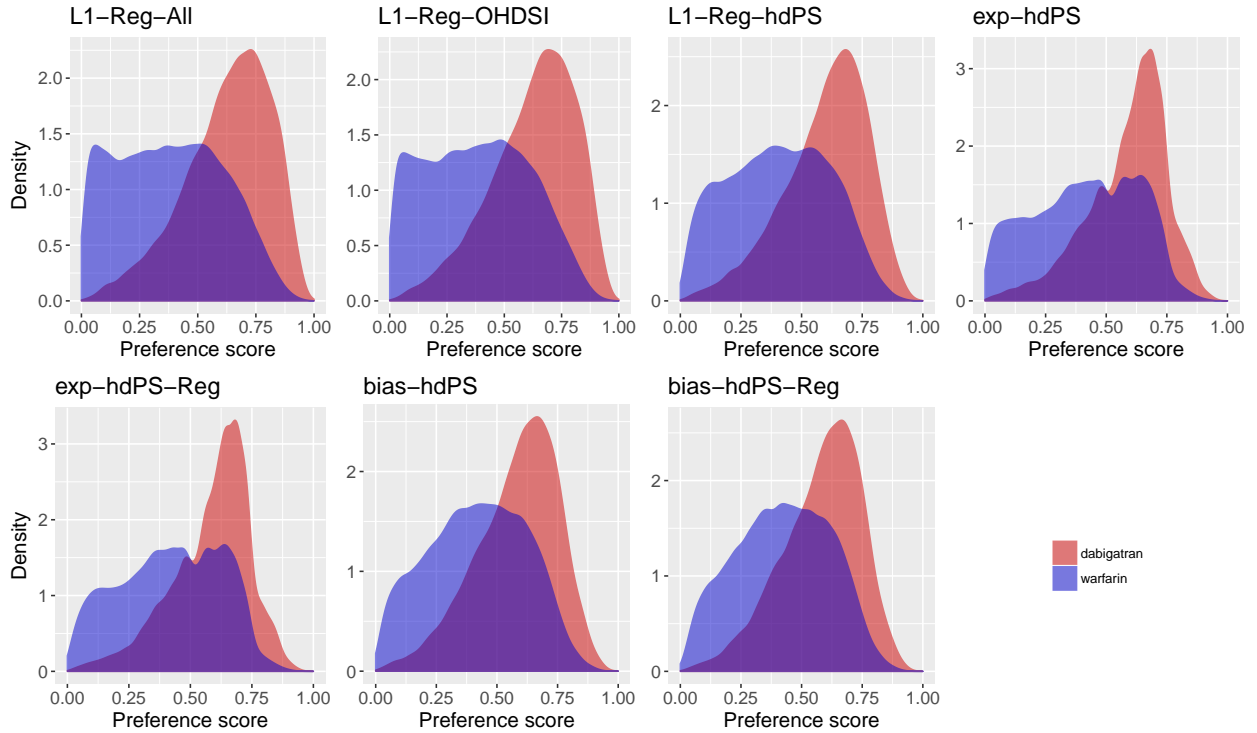
tributions and introduce different true effect sizes $\eta^*$. We use a decaying exponential function for both $S(t)$ and $C(t)$. We use a dataset of 50,000 assigned to 25,000 matched sets. For each effect size, we run 100 simulations.

| $\theta$ dist | $\eta^*$ | bias | sd |
|---|---|---|---|
| $\theta_{0,k}, \theta_{1,k} \sim N(0, 0.1)$ | -1.0 | 0.0026 | 0.0023 |
| $\theta_{0,k}, \theta_{1,k} \sim N(0, 0.1)$ | -0.5 | 0.0009 | 0.0019 |
| $\theta_{0,k}, \theta_{1,k} \sim N(0, 0.1)$ | 0.0 | -0.0002 | 0.0018 |
| $\theta_{0,k}, \theta_{1,k} \sim N(0, 0.1)$ | 0.5 | -0.0006 | 0.0018 |
| $\theta_{0,k}, \theta_{1,k} \sim N(0, 0.1)$ | 1.0 | -0.0045 | 0.0017 |
| $\theta_{0,k}, \theta_{1,k} \sim N(0, 1.0)$ | -1.0 | 0.1677 | 0.0022 |
| $\theta_{0,k}, \theta_{1,k} \sim N(0, 1.0)$ | -0.5 | 0.0910 | 0.0017 |
| $\theta_{0,k}, \theta_{1,k} \sim N(0, 1.0)$ | 0.0 | 0.000027 | 0.0017 |
| $\theta_{0,k}, \theta_{1,k} \sim N(0, 1.0)$ | 0.5 | -0.1007 | 0.0018 |
| $\theta_{0,k}, \theta_{1,k} \sim N(0, 1.0)$ | 1.0 | -0.2097 | 0.0015 |
| $\theta_{0,k}, \theta_{1,k} \sim N(1, 1.0)$ | -1.0 | 0.2085 | 0.0019 |
| $\theta_{0,k}, \theta_{1,k} \sim N(1, 1.0)$ | 0.0 | 0.0011 | 0.0017 |
| $\theta_{0,k}, \theta_{1,k} \sim N(1, 1.0)$ | 1.0 | -0.2380 | 0.0016 |
| $\theta_{0,k} \sim N(0, 1), \theta_{1,k} = \theta_{0,k}$ | -1.0 | 0.0035 | 0.0023 |
| $\theta_{0,k} \sim N(0, 1), \theta_{1,k} = \theta_{0,k}$ | 0.0 | -0.0006 | 0.0017 |
| $\theta_{0,k} \sim N(0, 1), \theta_{1,k} = \theta_{0,k}$ | 1.0 | 0.0023 | 0.0018 |

Supplementary Table 3: Estimation bias from true hazard ratio for different distributions of matched hazards

The simulations show that the observed bias arises from the distribution of hazard differences $\theta_{1,k} - \theta_{0,k}$ across matched sets. Even if the average hazard difference is exactly $\eta^*$, the variance in the hazard differences causes a bias towards the null. The larger the variance, the greater the bias.
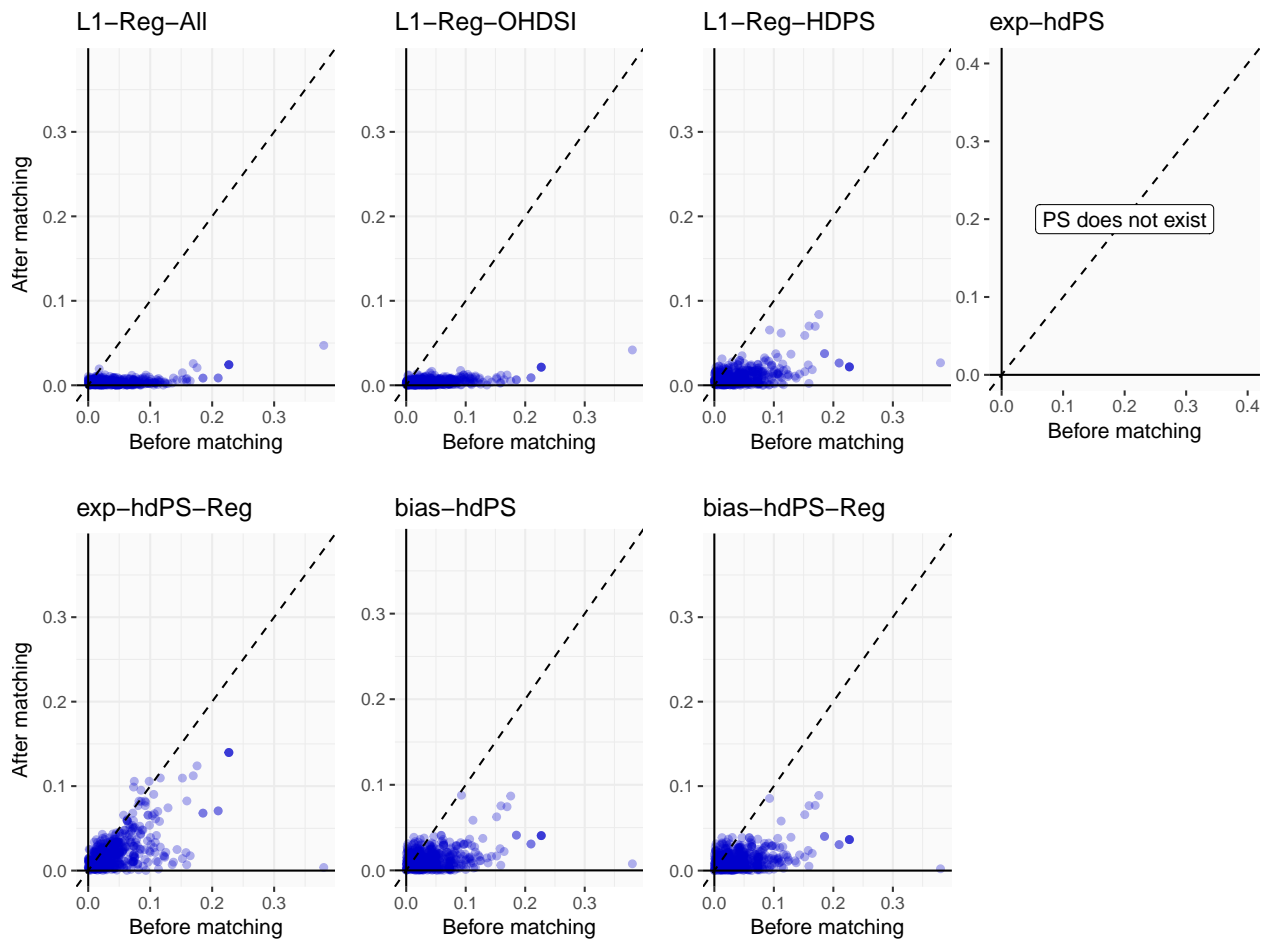
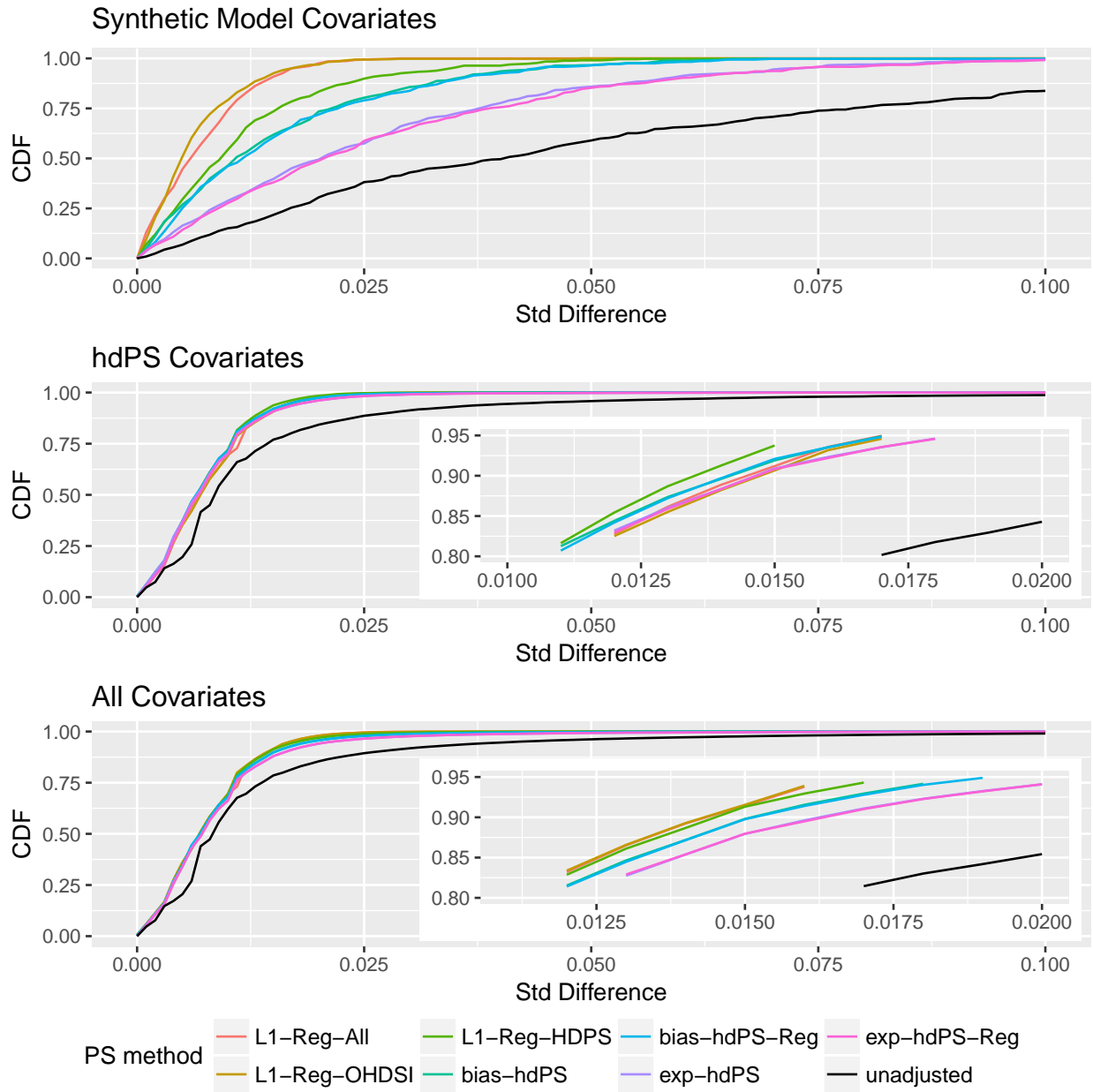## Supplementary Material 12: Additional Results

Supplementary Figure 1: Anticoagulants study: preference score distributions. Bias-based hdPS used on the empirical outcome of interest.
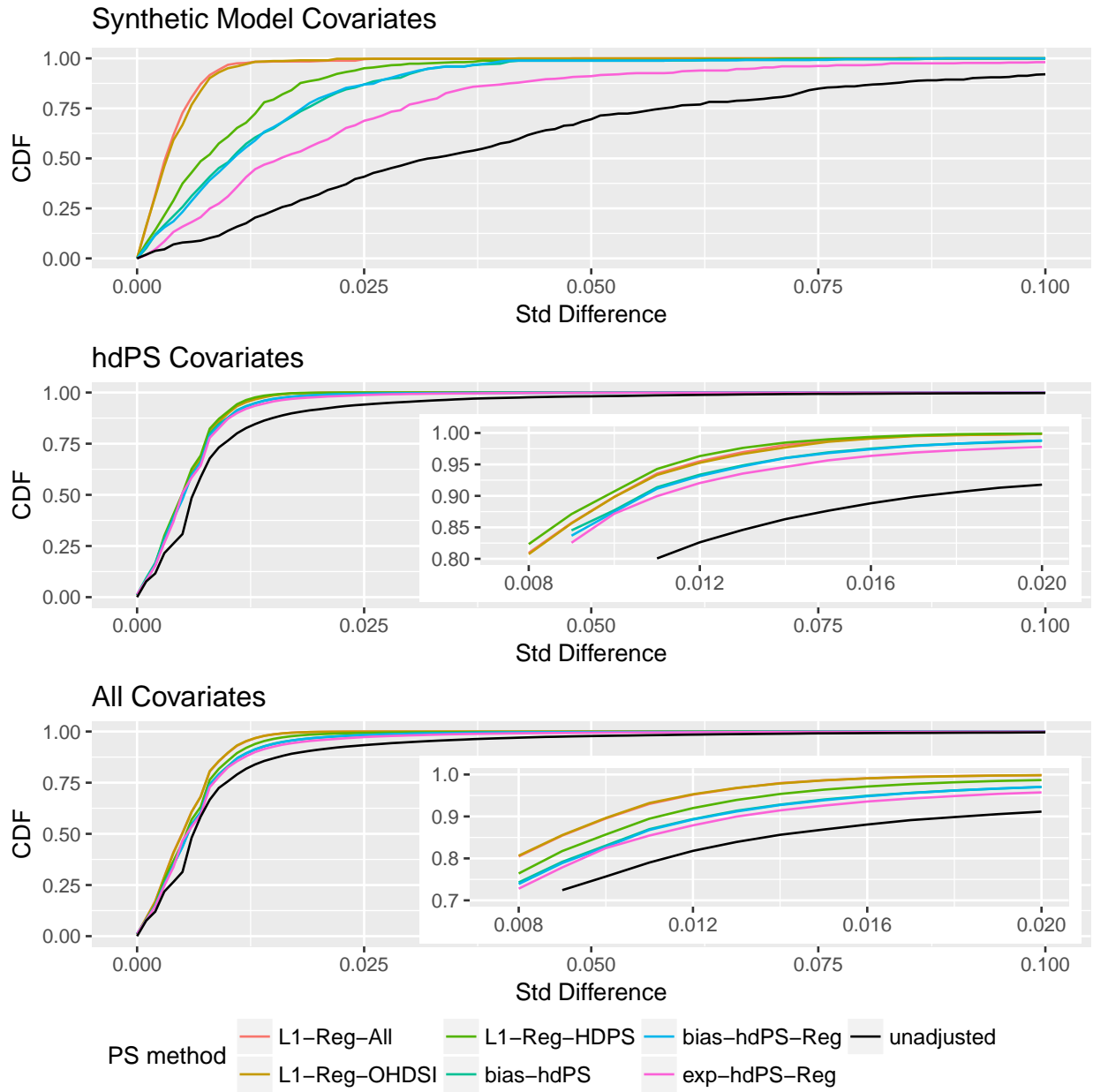


Supplementary Figure 2: NSAIDs study: preference score distributions. Bias-based hdPS used on the empirical outcome of interest.
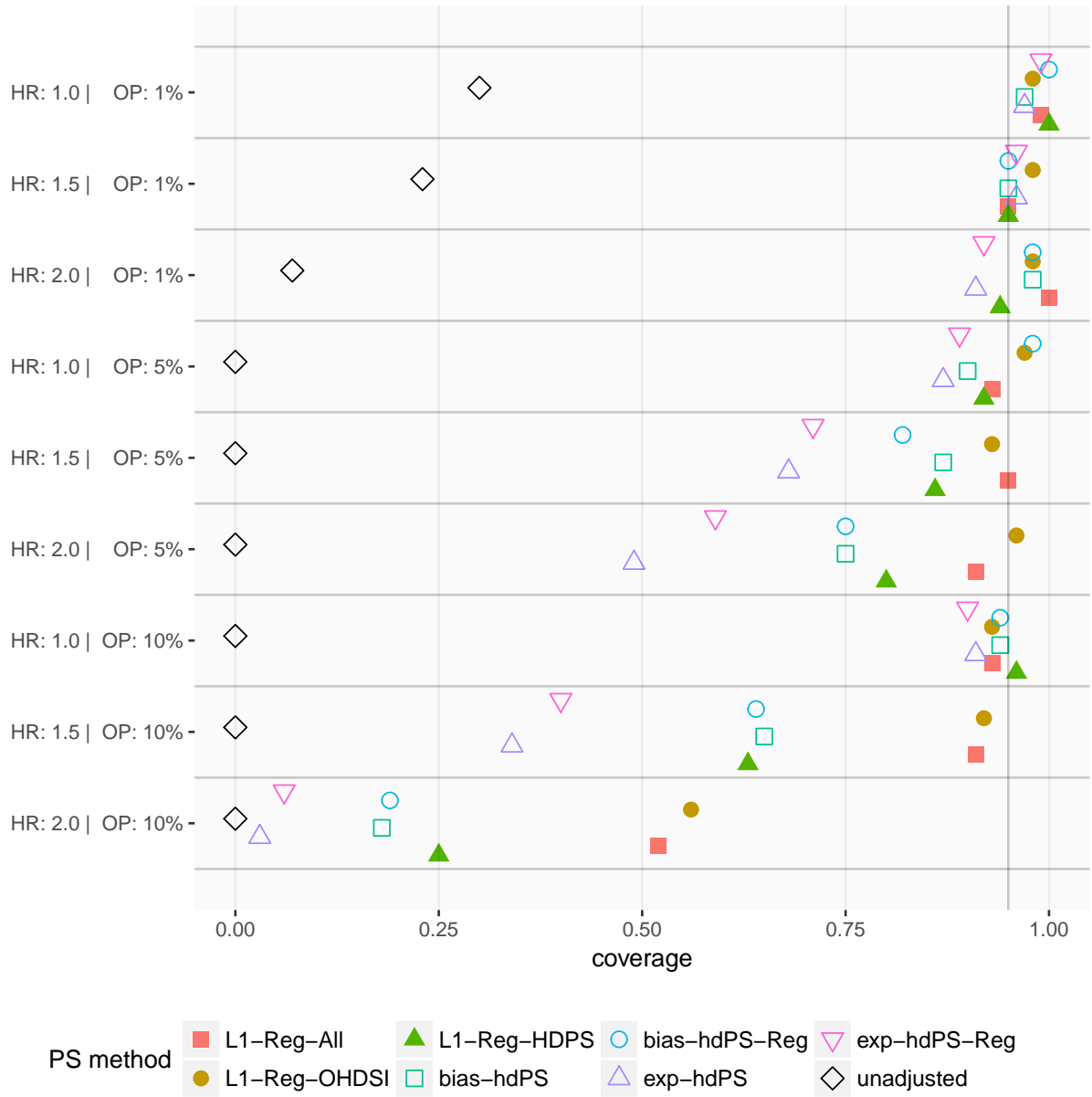
Supplementary Figure 3: NSAIDs study: before and after PS matching scatterplot of absolute standardized differences for synthetic model covariates. Before matching outlier is "Index Year: 2005," and corresponds to a documented decrease in celecoxib marketing and sales.

Supplementary Figure 4: Anticoagulants study: empirical cumulative distribution function of post-PS matching standardized differences for synthetic model covariates (top), hdPS Covariates (middle), and all covariates (bottom). Inlets provide closer view of differences between methods
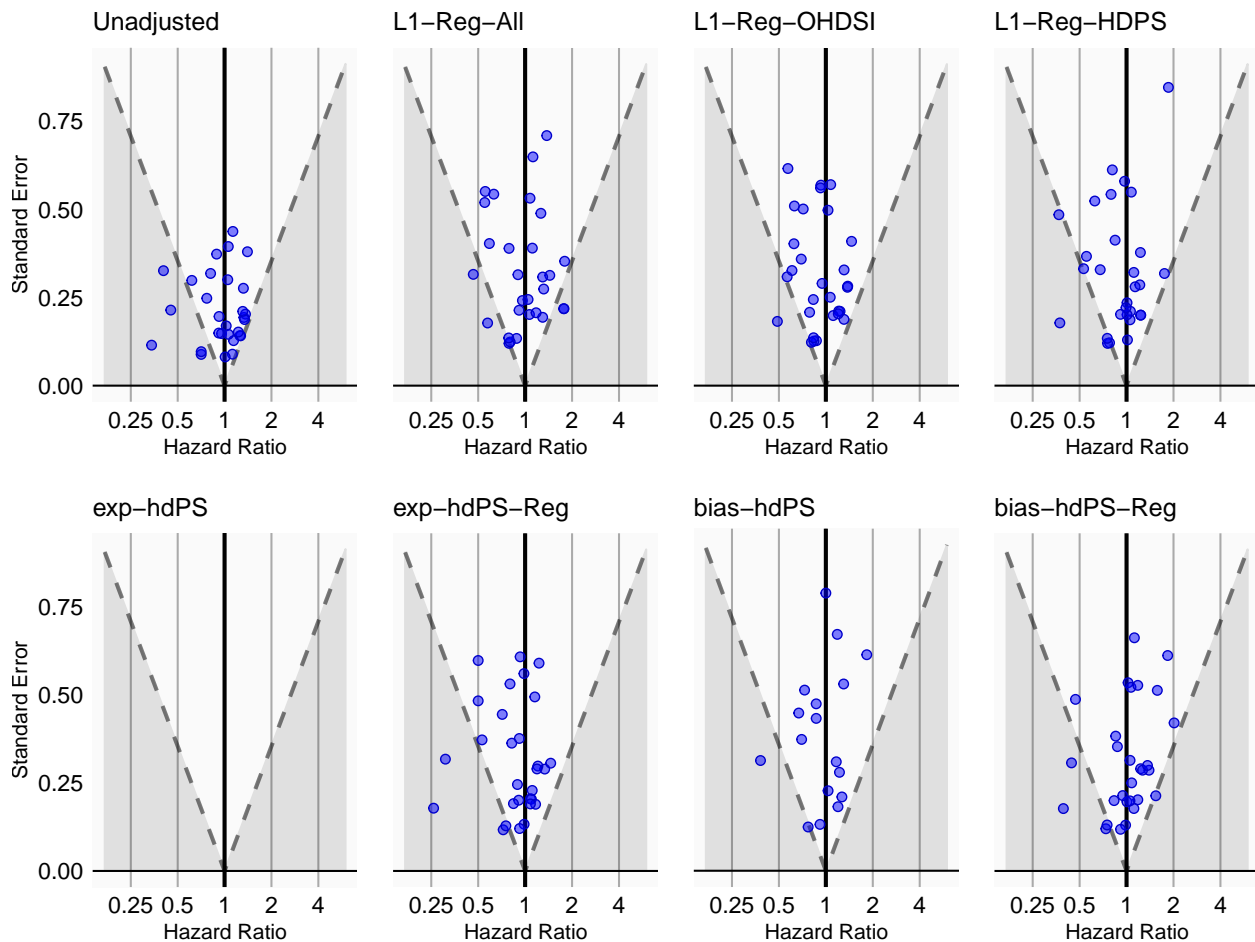
Supplementary Figure 5: NSAIDs study: empirical cumulative distribution function of post-PS matching standardized differences for synthetic model covariates (top), hdPS Covariates (middle), and all covariates (bottom). Inlets provide closer view of differences between methods

Supplementary Figure 6: Anticoagulants study: coverage of true hazard ratio (HR) across 100 simulations under different simulation parameters of true HR and outcome prevalence (OP)

Supplementary Figure 7: NSAIDs study: coverage of true hazard ratio (HR) across 100 simulations under different simulation parameters of true HR and outcome prevalence (OP)

Supplementary Figure 8: NSAIDs study: hazard ratio estimates (horizontal axis) and width of 95% confidence interval [via standard error] (vertical axis) for 29 negative control outcomes. Dashed line represents the straight line boundary at where the 95% confidence interval does (above) or does not (below) contain the assumed true hazard ratio of 1. Coverage indicates proportion of intervals that contains 1. bias-hdPS fails to construct 12 of the 29 PS models, and exp-hdPS fails to construct all 29.

# References

[1] Overhage JM, Ryan PB, Reich CG, Hartzema AG, and Stang PE (2011). Validation of a common data model for active safety surveillance research. *J Am Med Inform Assoc*, 19(1):54–60.

[2] Lawless J (1998). Parametric models in survival analysis. *Encyclopedia of Biostatistics*.

[3] Whittemore AS and Keller JB (1986). Survival estimation using splines. *Biometrics*, pages 495–506.

[4] Efron B (1977). The efficiency of Cox's likelihood function for censored data. *J Am Stat Assoc*, 72(359):557–565.

[5] Breslow N (1972). Discussion of the paper by D.R. Cox. *J R Stat Soc Series B Stat Methodol*, 34:216–217.

[6] Mittal S, Madigan D, Burd RS, and Suchard MA (2013). High-dimensional, massive sample-size Cox proportional hazards regression for survival analysis. *Biostatistics*, 15(2):207–221.

[7] Voss EA, Boyce RD, Ryan PB, van der Lei J, Rijnbeek PR, and Schuemie MJ (2017). Accuracy of an automated knowledge base for identifying drug adverse reactions. *J Biomed Inform*, 66:72–81.

[8] Avillach P, Dufour J.-C, Diallo G, et al. (2012). Design and validation of an automated method to detect known adverse drug reactions in MEDLINE: a contribution from the EU–ADR project. *Journal of the American Medical Informatics Association*, 20(3):446–452.

[9] Kilicoglu H, Rosemblat G, Fiszman M, and Rindflesch TC (2011). Constructing a semantic predication gold standard from the biomedical literature. *BMC bioinformatics*, 12(1):486.

[10] King G and Nielsen R (2015). Why propensity scores should not be used for matching. *Working Paper*.

[11] Austin PC and Stuart EA (2015). Optimal full matching for survival outcomes: a method that merits more widespread use. *Stat Med*, 34(30):3949–3967.

[12] Austin PC (2008). Assessing balance in measured baseline covariates when using many-to-one matching on the propensity-score. *Pharmacoepidemiol Drug Saf*, 17(12):1218–1225.

[13] Rassen JA, Shelat AA, Myers J, Glynn RJ, Rothman KJ, and Schneeweiss S (2012). One-to-many propensity score matching in cohort studies. *Pharmacoepidemiol Drug Saf*, 21(S2):69–80.