

# The hidden talents of logistic regression

From noisy labels to point processes

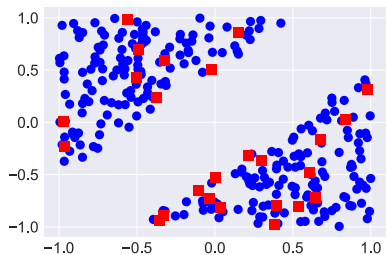
Aditya Krishna Menon



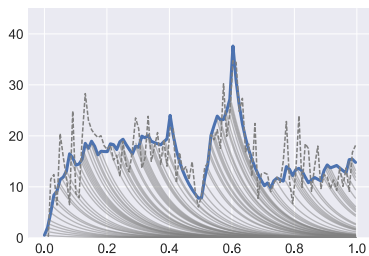
Australian  
National  
University

November 7th, 2017

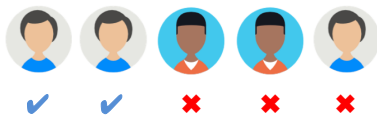
# Three problems...



*Label noise*

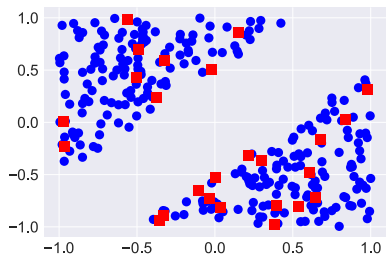


*Point processes*

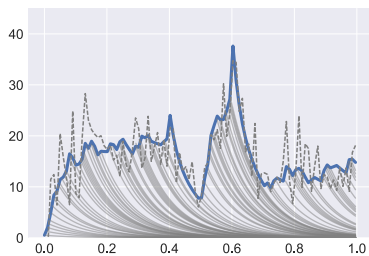


*Fairness*

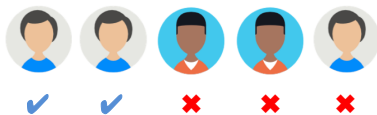
# Three problems...one solution?



*Label noise*

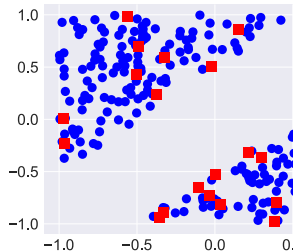


*Point processes*

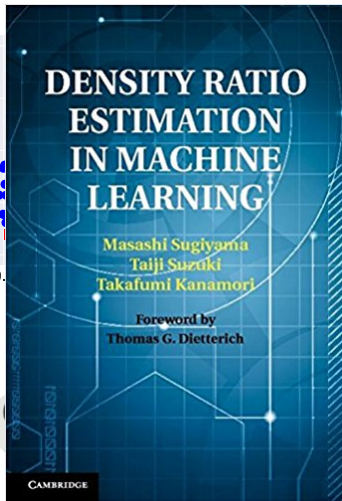


*Fairness*

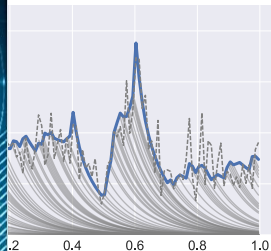
# Three problems...one solution?



*Label noise*



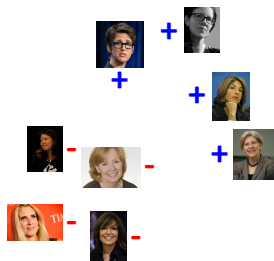
*Fairness*



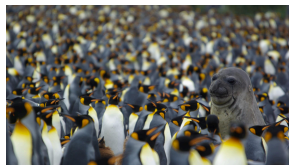
*Joint processes*



# DRE applications



*Covariate shift*

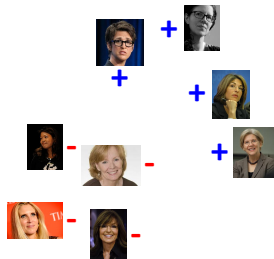


*Outlier detection*

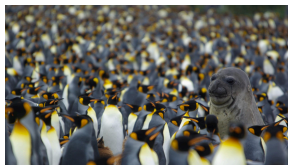


*Robot transition estimation*

# DRE applications



*Covariate shift*



*Outlier detection*

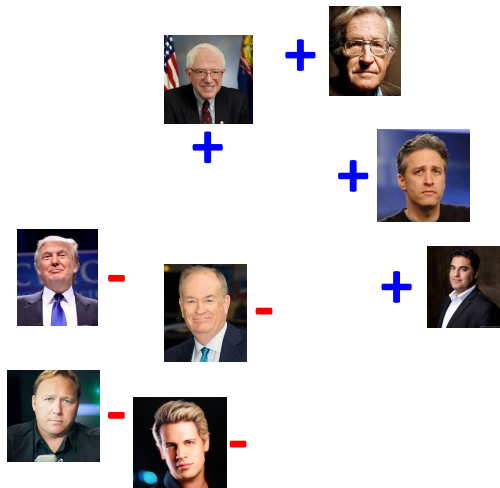


*Robot transition estimation*

In some cases, a different view may be more natural

# Class-probability estimation (CPE)

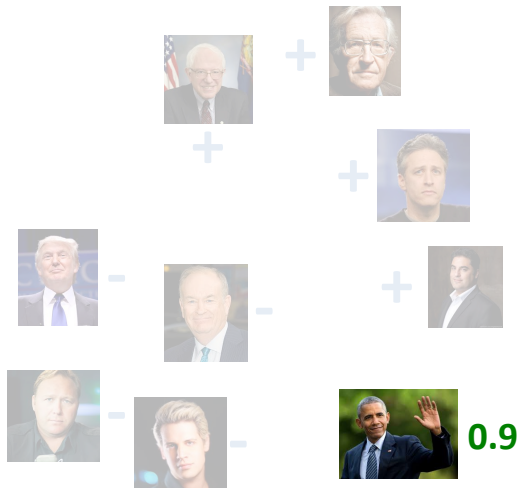
From labelled instances



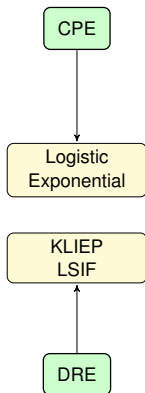
# Class-probability estimation (CPE)

From labelled instances, estimate **probability** of instance being + 've

- e.g. using logistic regression

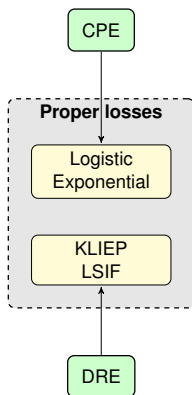


# This talk



# This talk

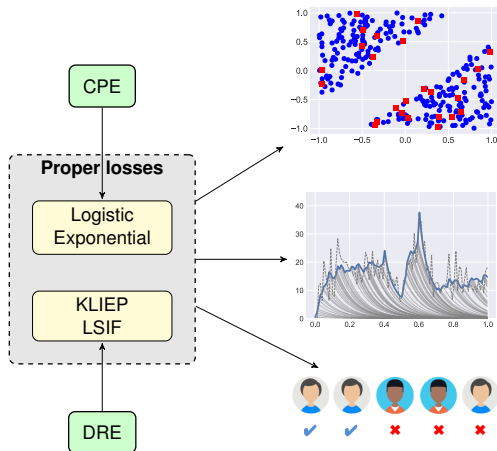
A formal link between DRE and CPE



# This talk

A formal link between DRE and CPE

CPE approach to three distinct learning problems



# Class-probability estimation



# Distributions for learning with binary labels

Fix an instance space  $\mathcal{X}$  (e.g.  $\mathbb{R}^n$ )

# Distributions for learning with binary labels

Fix an instance space  $\mathcal{X}$  (e.g.  $\mathbb{R}^n$ )

Let  $D$  be a distribution over  $\mathcal{X} \times \{\pm 1\}$ , with

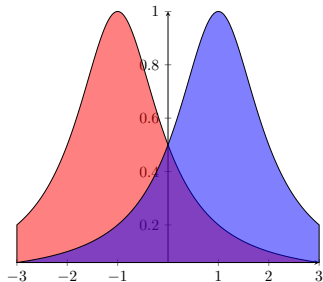
# Distributions for learning with binary labels

Fix an instance space  $\mathcal{X}$  (e.g.  $\mathbb{R}^n$ )

Let  $D$  be a distribution over  $\mathcal{X} \times \{\pm 1\}$ , with

$$(P(x), Q(x)) = (\mathbb{P}(\mathbf{X} = x \mid \mathbf{Y} = +1), \mathbb{P}(\mathbf{X} = x \mid \mathbf{Y} = -1))$$

Class conditionals



# Distributions for learning with binary labels

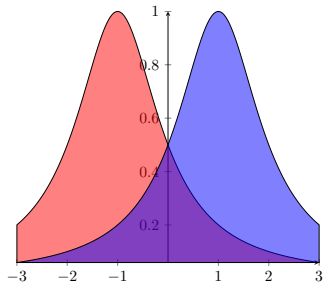
Fix an instance space  $\mathcal{X}$  (e.g.  $\mathbb{R}^n$ )

Let  $D$  be a distribution over  $\mathcal{X} \times \{\pm 1\}$ , with

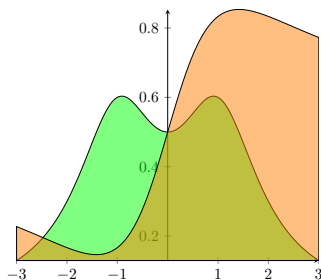
$$(P(x), Q(x)) = (\mathbb{P}(X = x \mid Y = +1), \mathbb{P}(X = x \mid Y = -1))$$

$$(M(x), \eta(x)) = (\mathbb{P}(X = x), \mathbb{P}(Y = +1 \mid X = x))$$

Class conditionals



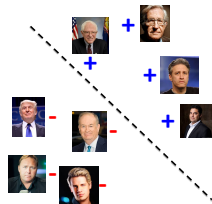
Marginal and class-probability function



# Scorers, losses, risks

A **scorer** is any  $s: \mathcal{X} \rightarrow \mathbb{R}$

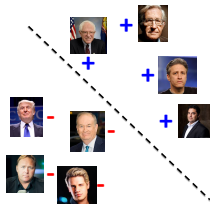
- e.g. linear scorer  $s: x \mapsto \langle w, x \rangle$



# Scorers, losses, risks

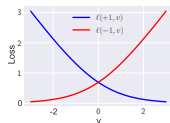
A **scorer** is any  $s: \mathcal{X} \rightarrow \mathbb{R}$

- e.g. linear scorer  $s: x \mapsto \langle w, x \rangle$



A **loss** is any  $\ell: \{\pm 1\} \times \mathbb{R} \rightarrow \mathbb{R}_+$

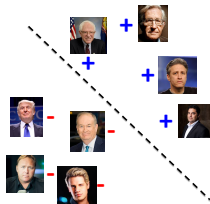
- e.g. logistic loss  $\ell: (y, v) \mapsto \log(1 + e^{-yv})$



# Scorers, losses, risks

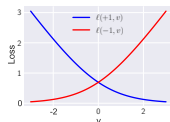
A **scorer** is any  $s: \mathcal{X} \rightarrow \mathbb{R}$

- e.g. linear scorer  $s: x \mapsto \langle w, x \rangle$



A **loss** is any  $\ell: \{\pm 1\} \times \mathbb{R} \rightarrow \mathbb{R}_+$

- e.g. logistic loss  $\ell: (y, v) \mapsto \log(1 + e^{-yv})$



The **risk** of scorer  $s$  wrt loss  $\ell$  and distribution  $D$  is

$$\mathbb{E}_{(X, Y) \sim D} [\ell(Y, s(X))]$$

- average loss on a random sample



# Class-probability estimation

**Goal:** estimate  $\eta(x) \doteq \mathbb{P}(Y = +1 \mid X = x)$



# Class-probability estimation

**Goal:** estimate  $\eta(x) \doteq \mathbb{P}(Y = +1 \mid X = x)$

For suitable  $\mathcal{S} \subset \mathbb{R}^{\mathcal{X}}$ , minimise empirical risk

$$\operatorname{argmin}_{s \in \mathcal{S}} \frac{1}{N} \sum_{n=1}^N \ell(y_n, s(x_n))$$

for strictly proper composite  $\ell$

# Class-probability estimation

**Goal:** estimate  $\eta(x) \doteq \mathbb{P}(Y = +1 \mid X = x)$

For suitable  $\mathcal{S} \subset \mathbb{R}^x$ , minimise empirical risk

$$\operatorname{argmin}_{s \in \mathcal{S}} \frac{1}{N} \sum_{n=1}^N \ell(y_n, s(x_n))$$

for strictly proper composite  $\ell$  viz., for invertible link  $\Psi : (0, 1) \rightarrow \mathbb{R}$ ,

$$\operatorname{argmin}_{s \in \mathbb{R}^x} \mathbb{E}[\ell(Y, s(X))] = \Psi \circ \eta$$

- e.g. for logistic loss,  $\Psi(u) = \log \frac{u}{1-u}$

# Class-probability estimation

**Goal:** estimate  $\eta(x) \doteq \mathbb{P}(Y = +1 \mid X = x)$

For suitable  $\mathcal{S} \subset \mathbb{R}^x$ , minimise empirical risk

$$\operatorname{argmin}_{s \in \mathcal{S}} \frac{1}{N} \sum_{n=1}^N \ell(y_n, s(x_n))$$

for strictly proper composite  $\ell$  viz., for invertible link  $\Psi : (0, 1) \rightarrow \mathbb{R}$ ,

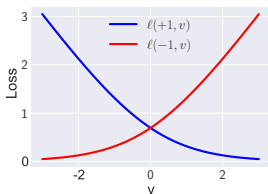
$$\operatorname{argmin}_{s \in \mathbb{R}^x} \mathbb{E}[\ell(Y, s(X))] = \Psi \circ \eta$$

- e.g. for logistic loss,  $\Psi(u) = \log \frac{u}{1-u}$

Estimate  $\hat{\eta} \doteq \Psi^{-1} \circ s$

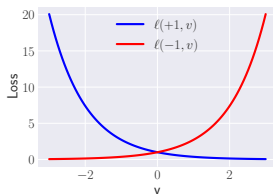
- e.g. for logistic loss,  $\hat{\eta}(x) = 1/(1 + \exp(-s(x)))$

# Examples of proper composite losses



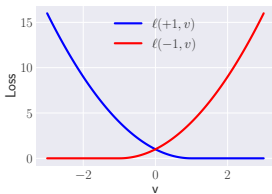
Logistic loss

$$\Psi^{-1} : v \mapsto 1/(1 + \exp(-v))$$



Exponential loss

$$\Psi^{-1} : v \mapsto 1/(1 + \exp(-2v))$$



Square hinge loss

$$\Psi^{-1} : v \mapsto \min(\max(0, (v + 1)/2), 1)$$

# Class-probabilities and density ratios

# CPE versus DRE

Given samples  $S \sim D^N$ , with  $D = (P, Q) = (M, \eta)$ :

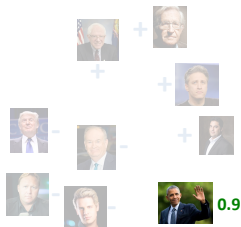
# CPE versus DRE

Given samples  $S \sim D^N$ , with  $D = (P, Q) = (M, \eta)$ :

## Class-probability estimation (CPE)

Estimate  $\eta$

- class-probability function



# CPE versus DRE

Given samples  $S \sim D^N$ , with  $D = (P, Q) = (M, \eta)$ :

## Class-probability estimation (CPE)

Estimate  $\eta$

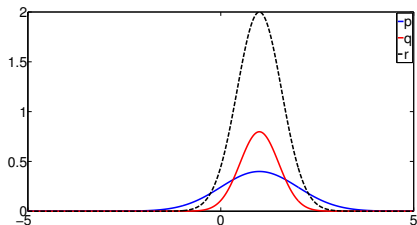
- class-probability function



## Density ratio estimation (DRE)

Estimate  $r = p/q$

- class-conditional density ratio





# CPE and DRE: exact solutions

Bayes' rule shows DRE and CPE are linked (Bickel et al, 2009):

$$r(x) \doteq \frac{p(x)}{q(x)}$$

# CPE and DRE: exact solutions

Bayes' rule shows DRE and CPE are linked (Bickel et al, 2009):

$$\begin{aligned} r(x) &\doteq \frac{p(x)}{q(x)} \\ &= \frac{\mathbb{P}(\mathbf{X} = x \mid \mathbf{Y} = +1)}{\mathbb{P}(\mathbf{X} = x \mid \mathbf{Y} = -1)} \end{aligned}$$

# CPE and DRE: exact solutions

Bayes' rule shows DRE and CPE are linked (Bickel et al, 2009):

$$\begin{aligned}r(x) &\doteq \frac{p(x)}{q(x)} \\ &= \frac{\mathbb{P}(\mathbf{X} = x \mid \mathbf{Y} = +1)}{\mathbb{P}(\mathbf{X} = x \mid \mathbf{Y} = -1)} \\ &\propto \frac{\mathbb{P}(\mathbf{Y} = +1 \mid \mathbf{X} = x)}{\mathbb{P}(\mathbf{Y} = -1 \mid \mathbf{X} = x)}\end{aligned}$$

# CPE and DRE: exact solutions

Bayes' rule shows DRE and CPE are linked (Bickel et al, 2009):

$$\begin{aligned}r(x) &\doteq \frac{p(x)}{q(x)} \\&= \frac{\mathbb{P}(\mathbf{X} = x \mid \mathbf{Y} = +1)}{\mathbb{P}(\mathbf{X} = x \mid \mathbf{Y} = -1)} \\&\propto \frac{\mathbb{P}(\mathbf{Y} = +1 \mid \mathbf{X} = x)}{\mathbb{P}(\mathbf{Y} = -1 \mid \mathbf{X} = x)} \\&= \frac{\eta(x)}{1 - \eta(x)}\end{aligned}$$

# CPE and DRE: exact solutions

Bayes' rule shows DRE and CPE are linked (Bickel et al, 2009):

$$\begin{aligned}r(x) &\doteq \frac{p(x)}{q(x)} \\&= \frac{\mathbb{P}(\mathbf{X} = x \mid \mathbf{Y} = +1)}{\mathbb{P}(\mathbf{X} = x \mid \mathbf{Y} = -1)} \\&\propto \frac{\mathbb{P}(\mathbf{Y} = +1 \mid \mathbf{X} = x)}{\mathbb{P}(\mathbf{Y} = -1 \mid \mathbf{X} = x)} \\&= \frac{\eta(x)}{1 - \eta(x)}\end{aligned}$$

Obtain  $\eta$  via CPE  $\rightarrow$  also obtain  $r$  for DRE

# CPE and DRE: exact solutions

Bayes' rule shows DRE and CPE are linked (Bickel et al, 2009):

$$\begin{aligned}r(x) &\doteq \frac{p(x)}{q(x)} \\&= \frac{\mathbb{P}(\mathbf{X} = x \mid \mathbf{Y} = +1)}{\mathbb{P}(\mathbf{X} = x \mid \mathbf{Y} = -1)} \\&\propto \frac{\mathbb{P}(\mathbf{Y} = +1 \mid \mathbf{X} = x)}{\mathbb{P}(\mathbf{Y} = -1 \mid \mathbf{X} = x)} \\&= \frac{\eta(x)}{1 - \eta(x)}\end{aligned}$$

Obtain  $\eta$  via CPE  $\rightarrow$  also obtain  $r$  for DRE

**But what about approximate solutions?**

# CPE and DRE: approximate solutions?

Natural class-probability estimate:  $\hat{\eta} \doteq \Psi^{-1} \circ s$

- e.g. for logistic loss,  $\hat{\eta}(x) = 1/(1 + e^{-s(x)})$

# CPE and DRE: approximate solutions?

Natural class-probability estimate:  $\hat{\eta} \doteq \Psi^{-1} \circ s$

- e.g. for logistic loss,  $\hat{\eta}(x) = 1/(1 + e^{-s(x)})$

Natural density ratio estimate:

$$\hat{r}(x) \doteq \frac{\hat{\eta}(x)}{1 - \hat{\eta}(x)}$$

- e.g. for logistic loss,  $\hat{r}(x) = e^{s(x)}$



# CPE and DRE: approximate solutions?

Natural class-probability estimate:  $\hat{\eta} \doteq \Psi^{-1} \circ s$

- e.g. for logistic loss,  $\hat{\eta}(x) = 1/(1 + e^{-s(x)})$

Natural density ratio estimate:

$$\hat{r}(x) \doteq \frac{\hat{\eta}(x)}{1 - \hat{\eta}(x)}$$

- e.g. for logistic loss,  $\hat{r}(x) = e^{s(x)}$

Intuitive, but what can we **guarantee** about this?

# CPE as Bregman minimisation

For proper composite  $\ell$ , the **regret** of a scorer is

$$\text{reg}(s) \doteq \mathbb{E}[\ell(\mathbf{Y}, s(\mathbf{X}))] - \min_{\bar{s} \in \mathbb{R}^{\mathcal{X}}} \mathbb{E}[\ell(\mathbf{Y}, \bar{s}(\mathbf{X}))]$$

# CPE as Bregman minimisation

For proper composite  $\ell$ , the **regret** of a scorer is

$$\begin{aligned}\text{reg}(s) &\doteq \mathbb{E}[\ell(\mathbf{Y}, s(\mathbf{X}))] - \min_{\bar{s} \in \mathbb{R}^{\mathcal{X}}} \mathbb{E}[\ell(\mathbf{Y}, \bar{s}(\mathbf{X}))] \\ &= \mathbb{E}[B_f(\boldsymbol{\eta}(\mathbf{X}), \hat{\boldsymbol{\eta}}(\mathbf{X}))]\end{aligned}$$

for **Bregman divergence**  $B_f(\cdot, \cdot)$  and loss-specific  $f$

# CPE as Bregman minimisation

For proper composite  $\ell$ , the **regret** of a scorer is

$$\begin{aligned}\text{reg}(s) &\doteq \mathbb{E}[\ell(\mathbf{Y}, s(\mathbf{X}))] - \min_{\bar{s} \in \mathbb{R}^{\mathcal{X}}} \mathbb{E}[\ell(\mathbf{Y}, \bar{s}(\mathbf{X}))] \\ &= \mathbb{E}[B_f(\boldsymbol{\eta}(\mathbf{X}), \hat{\boldsymbol{\eta}}(\mathbf{X}))]\end{aligned}$$

for **Bregman divergence**  $B_f(\cdot, \cdot)$  and loss-specific  $f$

- e.g. for logistic loss, regret is a KL projection

$$\text{reg}(s) = \mathbb{E}[\text{KL}(\boldsymbol{\eta}(\mathbf{X}) \parallel \hat{\boldsymbol{\eta}}(\mathbf{X}))]$$

# CPE as Bregman minimisation

For proper composite  $\ell$ , the **regret** of a scorer is

$$\begin{aligned}\text{reg}(s) &\doteq \mathbb{E}[\ell(Y, s(X))] - \min_{\bar{s} \in \mathbb{R}^x} \mathbb{E}[\ell(Y, \bar{s}(X))] \\ &= \mathbb{E}[B_f(\eta(X), \hat{\eta}(X))]\end{aligned}$$

for **Bregman divergence**  $B_f(\cdot, \cdot)$  and loss-specific  $f$

- e.g. for logistic loss, regret is a KL projection

$$\text{reg}(s) = \mathbb{E}[\text{KL}(\eta(X) \parallel \hat{\eta}(X))]$$

Concrete sense in which estimate  $\hat{\eta}$  is reasonable

# CPE as Bregman minimisation

For proper composite  $\ell$ , the **regret** of a scorer is

$$\begin{aligned}\text{reg}(s) &\doteq \mathbb{E}[\ell(\mathbf{Y}, s(\mathbf{X}))] - \min_{\bar{s} \in \mathbb{R}^x} \mathbb{E}[\ell(\mathbf{Y}, \bar{s}(\mathbf{X}))] \\ &= \mathbb{E}[B_f(\boldsymbol{\eta}(\mathbf{X}), \hat{\boldsymbol{\eta}}(\mathbf{X}))]\end{aligned}$$

for **Bregman divergence**  $B_f(\cdot, \cdot)$  and loss-specific  $f$

- e.g. for logistic loss, regret is a KL projection

$$\text{reg}(s) = \mathbb{E}[\text{KL}(\boldsymbol{\eta}(\mathbf{X}) \parallel \hat{\boldsymbol{\eta}}(\mathbf{X}))]$$

Concrete sense in which estimate  $\hat{\boldsymbol{\eta}}$  is reasonable

**Is there a similar sense in which  $\hat{\boldsymbol{\eta}}$  is reasonable?**

# A Bregman identity

One can show that:

$$(\forall x, y \in [0, 1]) B_f(x, y) =$$

# A Bregman identity

One can show that:

$$(\forall x, y \in [0, 1]) B_f(x, y) = (1 - x) \cdot B_{f^\oplus} \left( \frac{x}{1 - x}, \frac{y}{1 - y} \right),$$

where  $f^\oplus : z \mapsto (1 + z) \cdot f \left( \frac{z}{1 + z} \right)$



# A Bregman identity

One can show that:

$$(\forall x, y \in [0, 1]) B_f(x, y) = (1 - x) \cdot B_{f^\oplus} \left( \frac{x}{1 - x}, \frac{y}{1 - y} \right),$$

where  $f^\oplus : z \mapsto (1 + z) \cdot f \left( \frac{z}{1 + z} \right)$

Since  $r = \frac{\eta}{1 - \eta}$  and  $\hat{r} = \frac{\hat{\eta}}{1 - \hat{\eta}}$ ,

$$\text{reg}(s) = \mathbb{E}_{X \sim M} [B_f(\eta(X), \hat{\eta}(X))]$$

# A Bregman identity

One can show that:

$$(\forall x, y \in [0, 1]) B_f(x, y) = (1 - x) \cdot B_{f^\oplus} \left( \frac{x}{1 - x}, \frac{y}{1 - y} \right),$$

where  $f^\oplus : z \mapsto (1 + z) \cdot f \left( \frac{z}{1 + z} \right)$

Since  $r = \frac{\eta}{1 - \eta}$  and  $\hat{r} = \frac{\hat{\eta}}{1 - \hat{\eta}}$ ,

$$\text{reg}(s) = \mathbb{E}_{\mathbf{X} \sim M} [B_f(\eta(\mathbf{X}), \hat{\eta}(\mathbf{X}))] \propto \mathbb{E}_{\mathbf{X} \sim Q} [B_{f^\oplus}(r(\mathbf{X}), \hat{r}(\mathbf{X}))]$$

# A Bregman identity

One can show that:

$$(\forall x, y \in [0, 1]) B_f(x, y) = (1 - x) \cdot B_{f^\oplus} \left( \frac{x}{1 - x}, \frac{y}{1 - y} \right),$$

where  $f^\oplus : z \mapsto (1 + z) \cdot f \left( \frac{z}{1 + z} \right)$

Since  $r = \frac{\eta}{1 - \eta}$  and  $\hat{r} = \frac{\hat{\eta}}{1 - \hat{\eta}}$ ,

$$\text{reg}(s) = \mathbb{E}_{\mathbf{X} \sim M} [B_f(\eta(\mathbf{X}), \hat{\eta}(\mathbf{X}))] \propto \mathbb{E}_{\mathbf{X} \sim Q} [B_{f^\oplus}(r(\mathbf{X}), \hat{r}(\mathbf{X}))]$$

## CPE implicitly estimates density ratios

- complementary to (Sugiyama et al., 2012)

# CPE and DRE: summary

The asymptotic targets of CPE and DRE are closely linked

- $r = \frac{\eta}{1-\eta}$

# CPE and DRE: summary

The asymptotic targets of CPE and DRE are closely linked

- $r = \frac{\eta}{1-\eta}$

Estimating  $\hat{\eta}$  by proper loss minimisation is reasonable

- minimises Bregman divergence to  $\eta$

# CPE and DRE: summary

The asymptotic targets of CPE and DRE are closely linked

- $r = \frac{\eta}{1-\eta}$

Estimating  $\hat{\eta}$  by proper loss minimisation is reasonable

- minimises Bregman divergence to  $\eta$

Estimating  $\hat{r}$  by proper loss minimisation is reasonable

- minimises **alternate** Bregman divergence to  $r$

# CPE and DRE: summary

The asymptotic targets of CPE and DRE are closely linked

- $r = \frac{\eta}{1-\eta}$

Estimating  $\hat{\eta}$  by proper loss minimisation is reasonable

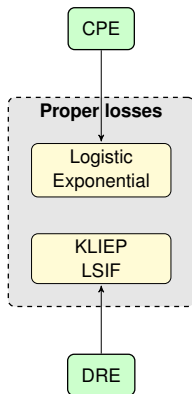
- minimises Bregman divergence to  $\eta$

Estimating  $\hat{r}$  by proper loss minimisation is reasonable

- minimises **alternate** Bregman divergence to  $r$

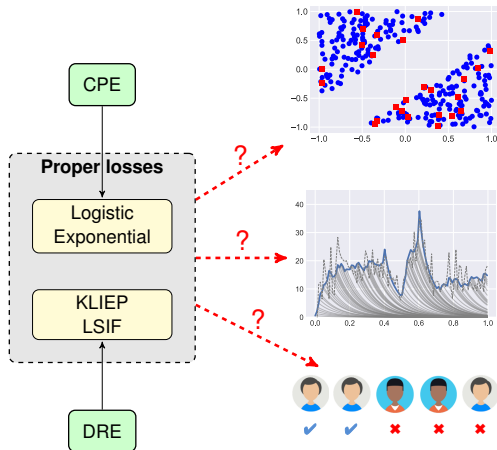
Underlying Bregman identity has multi-dimensional generalisation

# Summary thus far



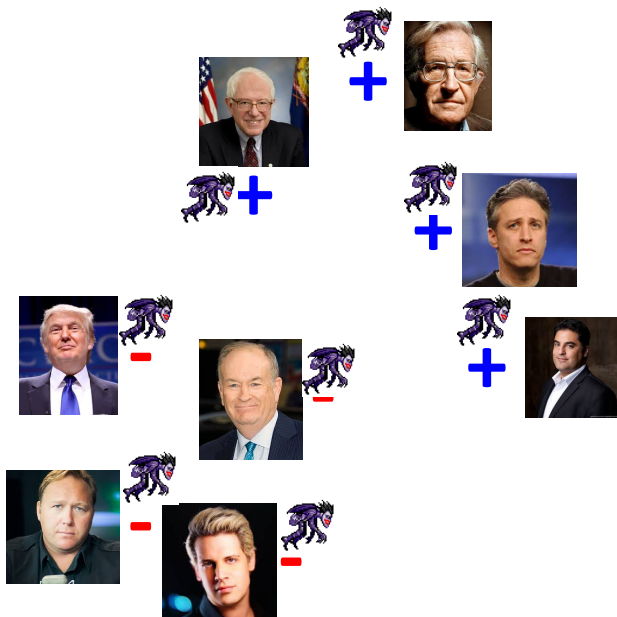


# Summary thus far

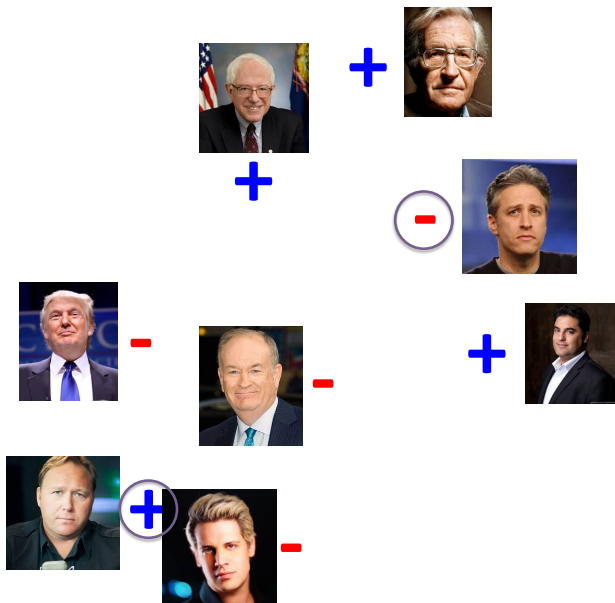


# Learning from noisy binary labels

# Learning from noisy binary labels



# Learning from noisy binary labels



# Label noise: formally

We care about “clean”  $D$

$$\min_s \mathbb{E}_{(X,Y)} [\ell(Y, s(X))]$$

**Ideal**

# Label noise: formally

We care about “clean”  $D$ , but observe samples from  $\bar{D} \neq D$

- $\mathbb{P}(X = x)$  remains same, but  $\mathbb{P}(Y = +1 | X = x)$  does not

**Ideal**

$$\min_s \mathbb{E}_{(X, Y)} [\ell(Y, s(X))]$$

**Reality**

$$\min_s \mathbb{E}_{(X, \bar{Y})} [\ell(\bar{Y}, s(X))]$$

# Label noise: formally

We care about “clean”  $D$ , but observe samples from  $\bar{D} \neq D$

- $\mathbb{P}(X = x)$  remains same, but  $\mathbb{P}(Y = +1 | X = x)$  does not

$$\min_s \mathbb{E}_{(X, Y)} [\ell(Y, s(X))]$$

**Ideal**

$$\min_s \mathbb{E}_{(X, \bar{Y})} [\ell(\bar{Y}, s(X))]$$

**Reality**

How to minimise the ideal risk?

# Label noise: a CPE perspective

Denote by  $\eta(x)$  the “clean”  $\mathbb{P}(Y = +1 \mid X = x)$



# Label noise: a CPE perspective

Denote by  $\eta(x)$  the “clean”  $\mathbb{P}(Y = +1 \mid X = x)$

Suppose  $(x, y)$  has label flipped with probability  $\rho \in [0, 1/2)$

## Label noise: a CPE perspective

Denote by  $\eta(x)$  the “clean”  $\mathbb{P}(Y = +1 \mid X = x)$

Suppose  $(x, y)$  has label flipped with probability  $\rho \in [0, 1/2)$

The “noisy”  $\mathbb{P}(\bar{Y} = +1 \mid X = x)$  is

$$\bar{\eta}(x) = (1 - \rho) \cdot \eta(x) + \rho \cdot (1 - \eta(x))$$

## Label noise: a CPE perspective

Denote by  $\eta(x)$  the “clean”  $\mathbb{P}(Y = +1 \mid X = x)$

Suppose  $(x, y)$  has label flipped with probability  $\rho \in [0, 1/2)$

The “noisy”  $\mathbb{P}(\bar{Y} = +1 \mid X = x)$  is

$$\bar{\eta}(x) = (1 - \rho) \cdot \eta(x) + \rho \cdot (1 - \eta(x))$$

We may write

$$\begin{bmatrix} \bar{\eta}(x) \\ 1 - \bar{\eta}(x) \end{bmatrix} = \underbrace{\begin{bmatrix} 1 - \rho & \rho \\ \rho & 1 - \rho \end{bmatrix}}_T \begin{bmatrix} \eta(x) \\ 1 - \eta(x) \end{bmatrix}$$

# Noise-corrected losses

$$\mathbb{E}_{(\mathbf{X}, \mathbf{Y})} [\ell(\mathbf{Y}, s(\mathbf{X}))]$$

## Noise-corrected losses

$$\mathbb{E}_{(\mathbf{X}, \mathbf{Y})} [\ell(\mathbf{Y}, s(\mathbf{X}))] = \mathbb{E}_{\mathbf{X} \sim M} \begin{bmatrix} \eta(\mathbf{X}) \\ 1 - \eta(\mathbf{X}) \end{bmatrix}^T [\ell(+1, s(\mathbf{X})) \quad \ell(-1, s(\mathbf{X}))]$$

## Noise-corrected losses

$$\begin{aligned}\mathbb{E}_{(\mathbf{X}, \mathbf{Y})} [\ell(\mathbf{Y}, s(\mathbf{X}))] &= \mathbb{E}_{\mathbf{X} \sim M} \begin{bmatrix} \eta(\mathbf{X}) \\ 1 - \eta(\mathbf{X}) \end{bmatrix}^T [\ell(+1, s(\mathbf{X})) \quad \ell(-1, s(\mathbf{X}))] \\ &= \mathbb{E}_{\mathbf{X} \sim M} \begin{bmatrix} \bar{\eta}(\mathbf{X}) \\ 1 - \bar{\eta}(\mathbf{X}) \end{bmatrix}^T (T^{-1})^T [\ell(+1, s(\mathbf{X})) \quad \ell(-1, s(\mathbf{X}))]\end{aligned}$$

## Noise-corrected losses

$$\begin{aligned}\mathbb{E}_{(\mathbf{X}, \mathbf{Y})} [\ell(\mathbf{Y}, s(\mathbf{X}))] &= \mathbb{E}_{\mathbf{X} \sim M} \begin{bmatrix} \eta(\mathbf{X}) \\ 1 - \eta(\mathbf{X}) \end{bmatrix}^T [\ell(+1, s(\mathbf{X})) \quad \ell(-1, s(\mathbf{X}))] \\ &= \mathbb{E}_{\mathbf{X} \sim M} \begin{bmatrix} \bar{\eta}(\mathbf{X}) \\ 1 - \bar{\eta}(\mathbf{X}) \end{bmatrix}^T (T^{-1})^T [\ell(+1, s(\mathbf{X})) \quad \ell(-1, s(\mathbf{X}))] \\ &= \mathbb{E}_{\mathbf{X} \sim M} \begin{bmatrix} \bar{\eta}(\mathbf{X}) \\ 1 - \bar{\eta}(\mathbf{X}) \end{bmatrix}^T [\bar{\ell}(+1, s(\mathbf{X})) \quad \bar{\ell}(-1, s(\mathbf{X}))]\end{aligned}$$

for the **noise-corrected loss** (Natarajan et al., 2013)

$$\bar{\ell}(y, v) = \frac{1}{1 - 2 \cdot \rho} ((1 - \rho) \cdot \ell(y, v) - \rho \cdot \ell(-y, v))$$

## Noise-corrected losses

$$\begin{aligned}\mathbb{E}_{(\mathbf{X}, \mathbf{Y})} [\ell(\mathbf{Y}, s(\mathbf{X}))] &= \mathbb{E}_{\mathbf{X} \sim M} \begin{bmatrix} \eta(\mathbf{X}) \\ 1 - \eta(\mathbf{X}) \end{bmatrix}^T [\ell(+1, s(\mathbf{X})) \quad \ell(-1, s(\mathbf{X}))] \\ &= \mathbb{E}_{\mathbf{X} \sim M} \begin{bmatrix} \bar{\eta}(\mathbf{X}) \\ 1 - \bar{\eta}(\mathbf{X}) \end{bmatrix}^T (T^{-1})^T [\ell(+1, s(\mathbf{X})) \quad \ell(-1, s(\mathbf{X}))] \\ &= \mathbb{E}_{\mathbf{X} \sim M} \begin{bmatrix} \bar{\eta}(\mathbf{X}) \\ 1 - \bar{\eta}(\mathbf{X}) \end{bmatrix}^T [\bar{\ell}(+1, s(\mathbf{X})) \quad \bar{\ell}(-1, s(\mathbf{X}))] \\ &= \mathbb{E}_{(\mathbf{X}, \bar{\mathbf{Y}})} [\bar{\ell}(\bar{\mathbf{Y}}, s(\mathbf{X}))],\end{aligned}$$

for the **noise-corrected loss** (Natarajan et al., 2013)

$$\bar{\ell}(y, v) = \frac{1}{1 - 2 \cdot \rho} ((1 - \rho) \cdot \ell(y, v) - \rho \cdot \ell(-y, v))$$



## Noise-corrected losses

$$\begin{aligned}\mathbb{E}_{(\mathbf{X}, \mathbf{Y})} [\ell(\mathbf{Y}, s(\mathbf{X}))] &= \mathbb{E}_{\mathbf{X} \sim M} \begin{bmatrix} \eta(\mathbf{X}) \\ 1 - \eta(\mathbf{X}) \end{bmatrix}^T [\ell(+1, s(\mathbf{X})) \quad \ell(-1, s(\mathbf{X}))] \\ &= \mathbb{E}_{\mathbf{X} \sim M} \begin{bmatrix} \bar{\eta}(\mathbf{X}) \\ 1 - \bar{\eta}(\mathbf{X}) \end{bmatrix}^T (T^{-1})^T [\ell(+1, s(\mathbf{X})) \quad \ell(-1, s(\mathbf{X}))] \\ &= \mathbb{E}_{\mathbf{X} \sim M} \begin{bmatrix} \bar{\eta}(\mathbf{X}) \\ 1 - \bar{\eta}(\mathbf{X}) \end{bmatrix}^T [\bar{\ell}(+1, s(\mathbf{X})) \quad \bar{\ell}(-1, s(\mathbf{X}))] \\ &= \mathbb{E}_{(\mathbf{X}, \bar{\mathbf{Y}})} [\bar{\ell}(\bar{\mathbf{Y}}, s(\mathbf{X}))],\end{aligned}$$

for the noise-corrected loss (Natarajan et al., 2013)

$$\bar{\ell}(y, v) = \frac{1}{1 - 2 \cdot \rho} ((1 - \rho) \cdot \ell(y, v) - \rho \cdot \ell(-y, v))$$

But  $\rho$  is unknown...

# Noise rate estimation

One can avoid knowing  $\rho$  for suitable  $\ell$

- eigenfunctions for the loss transform, e.g. “un-hinged” loss

# Noise rate estimation

One can avoid knowing  $\rho$  for suitable  $\ell$

- eigenfunctions for the loss transform, e.g. “un-hinged” loss

Alternately, assume  $\min_x \eta(x) = 0, \max_x \eta(x) = 1$

- “guaranteed” positive and negative instances
- c.f. ([Scott et al., 2013](#)), ([du Plessis et al., 2014](#))

# Noise rate estimation

One can avoid knowing  $\rho$  for suitable  $\ell$

- eigenfunctions for the loss transform, e.g. “un-hinged” loss

Alternately, assume  $\min_x \eta(x) = 0$ ,  $\max_x \eta(x) = 1$

- “guaranteed” positive and negative instances
- c.f. (Scott et al., 2013), (du Plessis et al., 2014)

Since  $\bar{\eta}(x) = (1 - 2 \cdot \rho) \cdot \eta(x) + \rho$ ,

$$\min_x \bar{\eta}(x) = \rho \quad \max_x \bar{\eta}(x) = 1 - \rho$$

# Noise rate estimation

One can avoid knowing  $\rho$  for suitable  $\ell$

- eigenfunctions for the loss transform, e.g. “un-hinged” loss

Alternately, assume  $\min_x \eta(x) = 0, \max_x \eta(x) = 1$

- “guaranteed” positive and negative instances
- c.f. (Scott et al., 2013), (du Plessis et al., 2014)

Since  $\bar{\eta}(x) = (1 - 2 \cdot \rho) \cdot \eta(x) + \rho$ ,

$$\min_x \bar{\eta}(x) = \rho \quad \max_x \bar{\eta}(x) = 1 - \rho$$

**Range of  $\bar{\eta}$  lets us estimate  $\rho$ !**

van Rooyen et al. Learning with symmetric label noise: the importance of being unhinged. NIPS 2015.

Menon et al. Learning from corrupted binary labels via class-probability estimation. ICML 2015.

# Beyond symmetric binary noise

For asymmetric multi-class noise, we similarly have

$$\bar{\eta}(x) = T\eta(x)$$

where e.g.  $\bar{\eta}(x) = (\mathbb{P}(Y = 1 | X = x), \dots, \mathbb{P}(Y = K | X = x))$

- analogous noise-corrected loss and noise estimation

# Beyond symmetric binary noise

For asymmetric multi-class noise, we similarly have

$$\bar{\eta}(x) = T\eta(x)$$

where e.g.  $\bar{\eta}(x) = (\mathbb{P}(Y = 1 | X = x), \dots, \mathbb{P}(Y = K | X = x))$

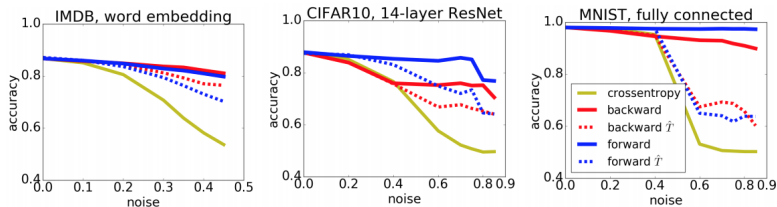
- analogous noise-corrected loss and noise estimation

Broader range of weakly supervised problems captured

- confer ([van Rooyen & Williamson, 2017](#))

# Illustration: deep network

## Corrected losses with and without noise estimation





## Instance-dependent noise?

Denote by  $\eta(x)$  the “clean”  $\mathbb{P}(Y = +1 \mid X = x)$

# Instance-dependent noise?

Denote by  $\eta(x)$  the “clean”  $\mathbb{P}(Y = +1 \mid X = x)$

Suppose  $(x, y)$  has label flipped with probability  $\rho(x) \in [0, 1/2)$

# Instance-dependent noise?

Denote by  $\eta(x)$  the “clean”  $\mathbb{P}(Y = +1 \mid X = x)$

Suppose  $(x, y)$  has label flipped with probability  $\rho(x) \in [0, 1/2)$

The “noisy” class-probability function is

$$\bar{\eta}(x) = (1 - \rho(x)) \cdot \eta(x) + \rho(x) \cdot (1 - \eta(x))$$

# Instance-dependent noise?

Denote by  $\eta(x)$  the “clean”  $\mathbb{P}(Y = +1 \mid X = x)$

Suppose  $(x, y)$  has label flipped with probability  $\rho(x) \in [0, 1/2)$

The “noisy” class-probability function is

$$\bar{\eta}(x) = (1 - \rho(x)) \cdot \eta(x) + \rho(x) \cdot (1 - \eta(x))$$

Estimating  $\rho(x)$  is non-trivial

- To make progress, we impose some structure on  $\rho$  and  $\eta$

# Assumptions on noise and distribution

Noise increases as  $\eta(x)$  approaches  $1/2$

- higher inherent uncertainty  $\rightarrow$  higher noise

# Assumptions on noise and distribution

Noise increases as  $\eta(x)$  approaches  $1/2$

- higher inherent uncertainty  $\rightarrow$  higher noise

Class-probability is expressible as

$$\eta(x) = u(\langle w^*, x \rangle)$$

for some non-decreasing, Lipschitz  $u(\cdot)$

- $u$  unknown  $\rightarrow$  single index model (SIM)
- such models learnable via **Isotron** (Kalai & Sastry, 2009)

# Structure of noisy class-probability

Under these assumptions, one may show

$$\bar{\eta}(x) = \bar{u}(\langle w^*, x \rangle)$$

for monotone  $\bar{u}$

- still in the SIM family!
- noise is baked into  $\bar{u}$

# Structure of noisy class-probability

Under these assumptions, one may show

$$\bar{\eta}(x) = \bar{u}(\langle w^*, x \rangle)$$

for monotone  $\bar{u}$

- still in the SIM family!
- noise is baked into  $\bar{u}$

One can estimate  $\bar{\eta}$  via Isotron

- do **not** need to know flip function  $\rho$  or link function  $u$



# Illustration: instance-dependent noise

Label flip function  $f(z) = (1 + e^{|z|/\alpha})^{-1}$

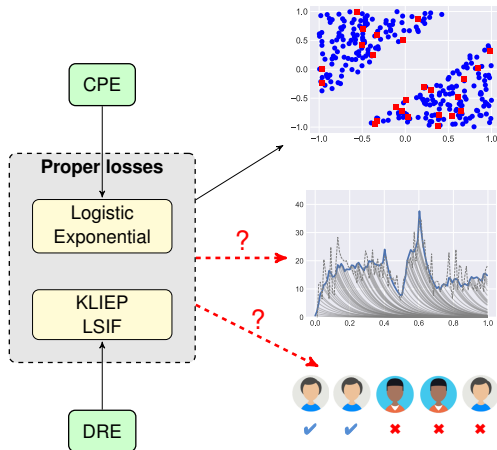
$\alpha$	Ridge ACC	Isotron ACC
$\frac{1}{8}$	0.9940 $\pm$ 0.0003	0.9974 $\pm$ 0.0002
$\frac{1}{4}$	0.9947 $\pm$ 0.0004	0.9974 $\pm$ 0.0003
$\frac{1}{2}$	0.9944 $\pm$ 0.0004	0.9937 $\pm$ 0.0006
1	0.9853 $\pm$ 0.0012	0.9700 $\pm$ 0.0021
2	0.8988 $\pm$ 0.0053	0.9239 $\pm$ 0.0050
4	0.7410 $\pm$ 0.0072	0.7863 $\pm$ 0.0138
8	0.6185 $\pm$ 0.0078	0.6467 $\pm$ 0.0405

usps 0v9

$\alpha$	Ridge ACC	Isotron ACC
$\frac{1}{8}$	0.9958 $\pm$ 0.0001	0.9984 $\pm$ 0.0001
$\frac{1}{4}$	0.9958 $\pm$ 0.0001	0.9979 $\pm$ 0.0001
$\frac{1}{2}$	0.9953 $\pm$ 0.0002	0.9966 $\pm$ 0.0003
1	0.9871 $\pm$ 0.0005	0.9864 $\pm$ 0.0007
2	0.9446 $\pm$ 0.0012	0.9565 $\pm$ 0.0013
4	0.8262 $\pm$ 0.0022	0.8768 $\pm$ 0.0041
8	0.6872 $\pm$ 0.0024	0.8088 $\pm$ 0.0291

mnist 6v7

# Summary thus far

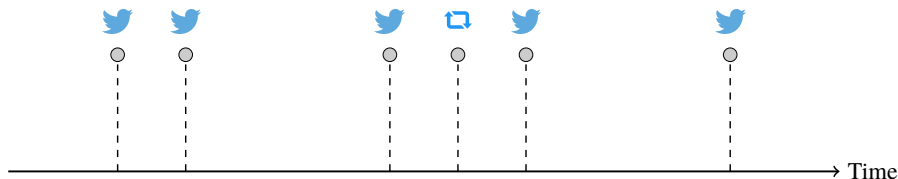


# Fitting point processes

# Point processes

Model the **rate** at which events occur in time

- re-tweets in a social network, earthquakes, ...



## Point processes: formally

Suppose  $(N(t))_{t \geq 0}$  counts the # of events in  $(0, t]$

# Point processes: formally

Suppose  $(N(t))_{t \geq 0}$  counts the # of events in  $(0, t]$

In the **non-homogeneous Poisson process (NHPP)**, one posits that the # of events in  $(s, t]$  follows

$$N(t) - N(s) \sim \text{Poiss} \left( \int_s^t \lambda(u) du \right)$$

for **intensity function**  $\lambda : \mathbb{R}_+ \rightarrow \mathbb{R}_+$

- instantaneous rate at which events occur

## NHPP likelihood

Suppose we observe event times  $\{t_1, \dots, t_N\}$ , with  $T \doteq \max_n t_n$

## NHPP likelihood

Suppose we observe event times  $\{t_1, \dots, t_N\}$ , with  $T \doteq \max_n t_n$

The negative log-likelihood for intensity  $\lambda(\cdot; \theta)$  is

$$\mathcal{L}(\theta) \doteq \sum_{n=1}^N -\log \lambda(t_n; \theta) + \int_0^T \lambda(u; \theta) \, du$$



## NHPP likelihood

Suppose we observe event times  $\{t_1, \dots, t_N\}$ , with  $T \doteq \max_n t_n$

The negative log-likelihood for intensity  $\lambda(\cdot; \theta)$  is

$$\begin{aligned}\mathcal{L}(\theta) &\doteq \sum_{n=1}^N -\log \lambda(t_n; \theta) + \int_0^T \lambda(u; \theta) \, du \\ &\propto \frac{1}{N} \sum_{n=1}^N -\log \lambda(t_n; \theta) + \frac{T}{N} \cdot \int_0^T \frac{1}{T} \cdot \lambda(u; \theta) \, du\end{aligned}$$

## NHPP likelihood

Suppose we observe event times  $\{t_1, \dots, t_N\}$ , with  $T \doteq \max_n t_n$

The negative log-likelihood for intensity  $\lambda(\cdot; \theta)$  is

$$\begin{aligned}\mathcal{L}(\theta) &\doteq \sum_{n=1}^N -\log \lambda(t_n; \theta) + \int_0^T \lambda(u; \theta) \, du \\ &\propto \frac{1}{N} \sum_{n=1}^N -\log \lambda(t_n; \theta) + \frac{T}{N} \cdot \int_0^T \frac{1}{T} \cdot \lambda(u; \theta) \, du \\ &= \mathbb{E}_{\mathbb{T} \sim \hat{P}} [-\log \lambda(\mathbb{T}; \theta)] + \frac{T}{N} \cdot \mathbb{E}_{\mathbb{T}' \sim Q} [\lambda(\mathbb{T}'; \theta)]\end{aligned}$$

where  $Q$  is uniform over  $[0, T]$

# NHPP likelihood

Suppose we observe event times  $\{t_1, \dots, t_N\}$ , with  $T \doteq \max_n t_n$

The negative log-likelihood for intensity  $\lambda(\cdot; \theta)$  is

$$\begin{aligned}\mathcal{L}(\theta) &\doteq \sum_{n=1}^N -\log \lambda(t_n; \theta) + \int_0^T \lambda(u; \theta) du \\ &\propto \frac{1}{N} \sum_{n=1}^N -\log \lambda(t_n; \theta) + \frac{T}{N} \cdot \int_0^T \frac{1}{T} \cdot \lambda(u; \theta) du \\ &= \mathbb{E}_{\mathbb{T} \sim \hat{P}} [-\log \lambda(\mathbb{T}; \theta)] + \frac{T}{N} \cdot \mathbb{E}_{\mathbb{T}' \sim Q} [\lambda(\mathbb{T}'; \theta)]\end{aligned}$$

where  $Q$  is uniform over  $[0, T]$

**Classification with a uniform background!**

# NHPPs as binary classification

On an interval  $[0, T]$ , event times  $\{t_1, \dots, t_N\}$  are iid with density

$$p(t) = \frac{\lambda(t)}{\int_0^T \lambda(u) du}$$

# NHPPs as binary classification

On an interval  $[0, T]$ , event times  $\{t_1, \dots, t_N\}$  are iid with density

$$p(t) = \frac{\lambda(t)}{\int_0^T \lambda(u) du}$$

Asymptotically, the likelihood is the classification risk

# NHPPs as binary classification

On an interval  $[0, T]$ , event times  $\{t_1, \dots, t_N\}$  are iid with density

$$p(t) = \frac{\lambda(t)}{\int_0^T \lambda(u) du}$$

Asymptotically, the likelihood is the classification risk with minimiser

$$\mathbb{E}_P[-\log \lambda(\mathbf{T})] + \frac{T}{N} \cdot \mathbb{E}_Q[\lambda(\mathbf{T}')] ]$$

# NHPPs as binary classification

On an interval  $[0, T]$ , event times  $\{t_1, \dots, t_N\}$  are iid with density

$$p(t) = \frac{\lambda(t)}{\int_0^T \lambda(u) du}$$

Asymptotically, the likelihood is the classification risk with minimiser

$$\operatorname{argmin}_{\lambda \in \mathbb{R}_+^{\mathcal{X}}} \mathbb{E}[-\log \lambda(\mathbf{T})] + \frac{T}{N} \cdot \mathbb{E}[\lambda(\mathbf{T}')] = N \cdot p$$

# NHPPs as binary classification

On an interval  $[0, T]$ , event times  $\{t_1, \dots, t_N\}$  are iid with density

$$p(t) = \frac{\lambda(t)}{\int_0^T \lambda(u) du}$$

Asymptotically, the likelihood is the classification risk with minimiser

$$\operatorname{argmin}_{\lambda \in \mathbb{R}_+^{\mathcal{X}}} \mathbb{E}[-\log \lambda(\mathbb{T})] + \frac{T}{N} \cdot \mathbb{E}[\lambda(\mathbb{T}')] = (N/T) \cdot p/q$$



# NHPPs as binary classification

On an interval  $[0, T]$ , event times  $\{t_1, \dots, t_N\}$  are iid with density

$$p(t) = \frac{\lambda(t)}{\int_0^T \lambda(u) du}$$

Asymptotically, the likelihood is the classification risk with minimiser

$$\operatorname{argmin}_{\lambda \in \mathbb{R}_+^{\mathcal{X}}} \mathbb{E}[-\log \lambda(\mathbf{T})] + \frac{T}{N} \cdot \frac{\mathbb{E}[\lambda(\mathbf{T}')] }{Q} = (N/T) \cdot p/q$$

Weighted density ratio estimation!

# Generalised likelihood?

For **scorer**  $s: \mathbb{R}_+ \rightarrow \mathbb{R}$ , consider

$$\begin{aligned} & \min_{s \in \mathcal{S}} \mathbb{E}_{\hat{P}} [\ell(+1, s(\mathbb{T}))] + \frac{T}{N} \cdot \mathbb{E}_{\mathcal{Q}} [\ell(-1, s(\mathbb{T}'))] \\ &= \min_{s \in \mathcal{S}} \sum_{n=1}^N \ell(+1, s(t_n)) + \int_0^T \ell(-1, s(t)) dt \end{aligned}$$

for strictly proper composite  $\ell$

# Generalised likelihood?

For **scorer**  $s: \mathbb{R}_+ \rightarrow \mathbb{R}$ , consider

$$\begin{aligned} & \min_{s \in \mathcal{S}} \mathbb{E}[\ell(+1, s(\mathbb{T}))] + \frac{T}{N} \cdot \mathbb{E}[\ell(-1, s(\mathbb{T}'))] \\ &= \min_{s \in \mathcal{S}} \sum_{n=1}^N \ell(+1, s(t_n)) + \int_0^T \ell(-1, s(t)) dt \end{aligned}$$

for strictly proper composite  $\ell$

We retain the optimal solution by picking

$$\lambda(t) = \frac{\Psi^{-1}(s(t))}{1 - \Psi^{-1}(s(t))}$$

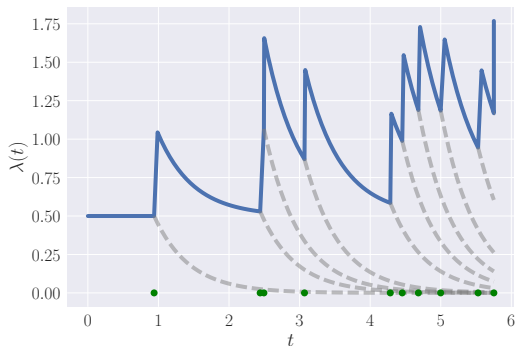
- optimal  $s = \Psi(\eta)$ ,  $\frac{\eta}{1-\eta} \propto p$

## Application: Hawkes processes

The **self-exciting** Hawkes process assumes, for link  $F(\cdot)$ ,

$$\lambda(t; \{t_n\}_{n=1}^N) = F\left(\mu + \alpha \cdot \sum_{t_n < t} e^{-\delta \cdot (t-t_n)}\right)$$

- occurrence of one event triggers subsequent events



# Generalised Hawkes likelihood?

In terms of a scorer, the Hawkes intensity is

$$\begin{aligned}\lambda(t; \{t_n\}_{n=1}^N) &= F(s(t)) \\ s(t) &= \mu + \alpha \cdot \Phi(t) \\ \Phi(t) &\doteq \sum_{t_n < t} e^{-\delta \cdot (t - t_n)}\end{aligned}$$

# Generalised Hawkes likelihood?

In terms of a scorer, the Hawkes intensity is

$$\begin{aligned}\lambda(t; \{t_n\}_{n=1}^N) &= F(s(t)) \\ s(t) &= \mu + \alpha \cdot \Phi(t) \\ \Phi(t) &\doteq \sum_{t_n < t} e^{-\delta \cdot (t - t_n)}\end{aligned}$$

Can minimise a proper loss with this  $s(\cdot)$  and  $\Phi$ , and set

$$\lambda(t) = \frac{\Psi^{-1}(s(t))}{1 - \Psi^{-1}(s(t))}$$

# Generalised Hawkes likelihood?

In terms of a scorer, the Hawkes intensity is

$$\begin{aligned}\lambda(t; \{t_n\}_{n=1}^N) &= F(s(t)) \\ s(t) &= \mu + \alpha \cdot \Phi(t) \\ \Phi(t) &\doteq \sum_{t_n < t} e^{-\delta \cdot (t - t_n)}\end{aligned}$$

Can minimise a proper loss with this  $s(\cdot)$  and  $\Phi$ , and set

$$\lambda(t) = \frac{\Psi^{-1}(s(t))}{1 - \Psi^{-1}(s(t))} = F(s(t))$$

if we choose

$$\Psi^{-1}(v) = \frac{F(v)}{1 + F(v)}$$

## Hawkes process with linear $F(\cdot)$

For  $F(z) = z$ , we may explore losses with  $\Psi(u) = \frac{u}{1-u}$

- losses that directly seek density ratios



## Hawkes process with linear $F(\cdot)$

For  $F(z) = z$ , we may explore losses with  $\Psi(u) = \frac{u}{1-u}$

- losses that directly seek density ratios

One appealing candidate ([Kanamori et al., 2009](#)):

$$\ell(+1, v) = -v \quad \ell(-1, v) = \frac{1}{2}v^2$$

- c.f. ([Reynaud-Bouret 2014](#), [Bacry et al., 2015](#))

## Hawkes process with linear $F(\cdot)$

For  $F(z) = z$ , we may explore losses with  $\Psi(u) = \frac{u}{1-u}$

- losses that directly seek density ratios

One appealing candidate ([Kanamori et al., 2009](#)):

$$\ell(+1, v) = -v \quad \ell(-1, v) = \frac{1}{2}v^2$$

- c.f. ([Reynaud-Bouret 2014](#), [Bacry et al., 2015](#))

Potential **closed-form** solution

$$\theta^* = \frac{N}{T} \cdot \left( \mathbb{E}_{\hat{Q}} [\Phi(\mathbb{T}') \Phi(\mathbb{T}')^T] \right)^{-1} \mathbb{E}_{\hat{P}} [\Phi(\mathbb{T})]$$

when this quantity is non-negative

## Hawkes process with exponential $F(\cdot)$

For  $F(z) = e^z$ , we may explore losses with  $\Psi(u) = \log \frac{u}{1-u}$

# Hawkes process with exponential $F(\cdot)$

For  $F(z) = e^z$ , we may explore losses with  $\Psi(u) = \log \frac{u}{1-u}$

One appealing candidate is familiar **logistic loss**

- nonlinear Hawkes with logistic regression!

# Hawkes process with exponential $F(\cdot)$

For  $F(z) = e^z$ , we may explore losses with  $\Psi(u) = \log \frac{u}{1-u}$

One appealing candidate is familiar **logistic loss**

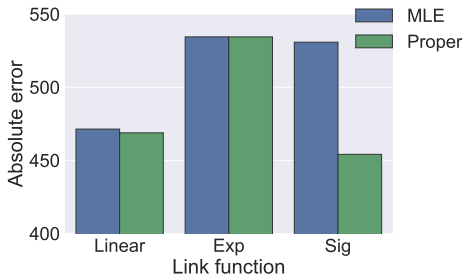
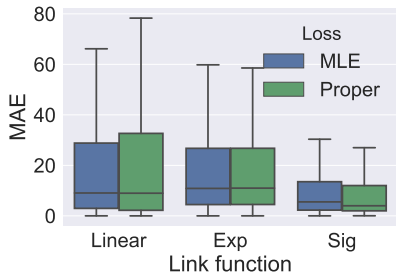
- nonlinear Hawkes with logistic regression!

By **weighting** the negative class, this is actually **equivalent** to MLE

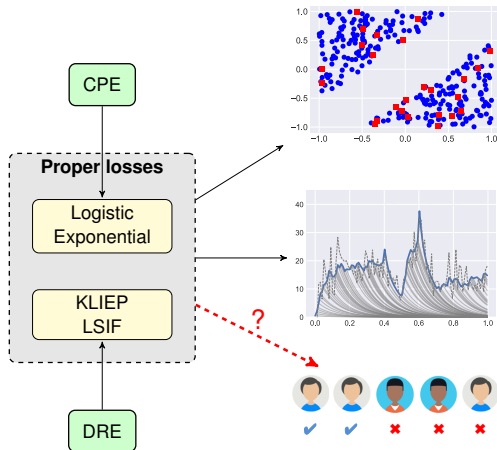
- follows from ([Fithian & Hastie, 2013](#))

# Illustration: fitting with proper losses

Prediction of # events on lastfm and bitcoin datasets



# Summary thus far



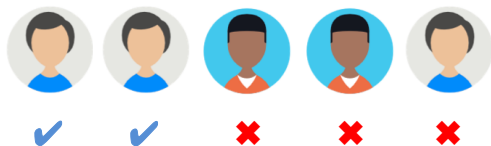
# Fairness-aware classification



# Fairness-aware classification

Learn a classifier achieving two goals:

- accurately predict a target label
- don't **discriminate** on some sensitive feature



# Fairness-aware classification: formally

We seek a classifier  $f: \mathcal{X} \rightarrow \{\pm 1\}$ , with induced predictions  $\hat{Y}$

# Fairness-aware classification: formally

We seek a classifier  $f: \mathcal{X} \rightarrow \{\pm 1\}$ , with induced predictions  $\hat{Y}$

$f$  should predict **well** target variable  $Y$

- e.g. attain **low** balanced error,

$$\text{BER}(f) \doteq \frac{1}{2} \cdot (\mathbb{P}(\hat{Y} = +1 \mid Y = -1) + \mathbb{P}(\hat{Y} = -1 \mid Y = +1))$$

# Fairness-aware classification: formally

We seek a classifier  $f: \mathcal{X} \rightarrow \{\pm 1\}$ , with induced predictions  $\hat{Y}$

$f$  should predict **well** target variable  $Y$

- e.g. attain **low** balanced error,

$$\text{BER}(f) \doteq \frac{1}{2} \cdot (\mathbb{P}(\hat{Y} = +1 \mid Y = -1) + \mathbb{P}(\hat{Y} = -1 \mid Y = +1))$$

$f$  should predict **poorly** sensitive variable  $\bar{Y}$

- e.g. attain **high** balanced error,

$$\overline{\text{BER}}(f) \doteq \frac{1}{2} \cdot (\mathbb{P}(\hat{Y} = +1 \mid \bar{Y} = -1) + \mathbb{P}(\hat{Y} = -1 \mid \bar{Y} = +1))$$

# Fairness-aware objective

We seek a solution to

$$\min_f \text{BER}(f) - \lambda \cdot \overline{\text{BER}}(f)$$

# Fairness-aware objective

We seek a solution to

$$\begin{aligned} \min_f \text{BER}(f) - \lambda \cdot \overline{\text{BER}}(f) &= \min_s \mathbb{E}_{\frac{P}{\lambda}} \llbracket s(\mathbf{X}) < 0 \rrbracket + \mathbb{E}_{\frac{Q}{\lambda}} \llbracket s(\mathbf{X}') > 0 \rrbracket \\ &\quad - \lambda \cdot \left( \mathbb{E}_{\frac{P}{\lambda}} \llbracket s(\mathbf{X}) < 0 \rrbracket + \mathbb{E}_{\frac{Q}{\lambda}} \llbracket s(\mathbf{X}') > 0 \rrbracket \right) \end{aligned}$$

# Fairness-aware objective

We seek a solution to

$$\begin{aligned} \min_f \text{BER}(f) - \lambda \cdot \overline{\text{BER}}(f) &= \min_s \mathbb{E}_{\underline{P}} \llbracket s(\mathbf{X}) < 0 \rrbracket + \mathbb{E}_{\underline{Q}} \llbracket s(\mathbf{X}') > 0 \rrbracket \\ &\quad - \lambda \cdot \left( \mathbb{E}_{\underline{P}} \llbracket s(\mathbf{X}) < 0 \rrbracket + \mathbb{E}_{\underline{Q}} \llbracket s(\mathbf{X}') > 0 \rrbracket \right) \end{aligned}$$

Natural to consider surrogate risk

$$\min_s \text{BER}_\ell(s) - \lambda \cdot \overline{\text{BER}}_\ell(s)$$

# Fairness-aware objective

We seek a solution to

$$\begin{aligned} \min_f \text{BER}(f) - \lambda \cdot \overline{\text{BER}}(f) &= \min_s \mathbb{E}_{\underline{P}} \llbracket s(\mathbf{X}) < 0 \rrbracket + \mathbb{E}_{\underline{Q}} \llbracket s(\mathbf{X}') > 0 \rrbracket \\ &\quad - \lambda \cdot \left( \mathbb{E}_{\underline{P}} \llbracket s(\mathbf{X}) < 0 \rrbracket + \mathbb{E}_{\underline{Q}} \llbracket s(\mathbf{X}') > 0 \rrbracket \right) \end{aligned}$$

Natural to consider surrogate risk

$$\begin{aligned} \min_s \text{BER}_\ell(s) - \lambda \cdot \overline{\text{BER}}_\ell(s) &= \min_s \mathbb{E}_{\underline{P}} \ell(+1, s(\mathbf{X})) + \mathbb{E}_{\underline{Q}} \ell(-1, s(\mathbf{X}')) \\ &\quad - \lambda \cdot \left( \mathbb{E}_{\underline{P}} \ell(+1, s(\mathbf{X})) + \mathbb{E}_{\underline{Q}} \ell(-1, s(\mathbf{X}')) \right) \end{aligned}$$



# Fairness-aware objective

We seek a solution to

$$\begin{aligned} \min_f \text{BER}(f) - \lambda \cdot \overline{\text{BER}}(f) &= \min_s \mathbb{E}_{\overline{P}} \llbracket s(\mathbf{X}) < 0 \rrbracket + \mathbb{E}_{\overline{Q}} \llbracket s(\mathbf{X}') > 0 \rrbracket \\ &\quad - \lambda \cdot \left( \mathbb{E}_{\overline{P}} \llbracket s(\mathbf{X}) < 0 \rrbracket + \mathbb{E}_{\overline{Q}} \llbracket s(\mathbf{X}') > 0 \rrbracket \right) \end{aligned}$$

Natural to consider surrogate risk

$$\begin{aligned} \min_s \text{BER}_\ell(s) - \lambda \cdot \overline{\text{BER}}_\ell(s) &= \min_s \mathbb{E}_{\overline{P}} \ell(+1, s(\mathbf{X})) + \mathbb{E}_{\overline{Q}} \ell(-1, s(\mathbf{X}')) \\ &\quad - \lambda \cdot \left( \mathbb{E}_{\overline{P}} \ell(+1, s(\mathbf{X})) + \mathbb{E}_{\overline{Q}} \ell(-1, s(\mathbf{X}')) \right) \end{aligned}$$

but in general this will be non-convex

# CPE approach?

Alternately, let us consider the Bayes-optimal solutions

$$f^* \in \operatorname{argmin}_{f: \mathcal{X} \rightarrow \{\pm 1\}} \operatorname{BER}(f) - \lambda \cdot \overline{\operatorname{BER}}(f)$$

# CPE approach?

Alternately, let us consider the Bayes-optimal solutions

$$f^* \in \operatorname{argmin}_{f: \mathcal{X} \rightarrow \{\pm 1\}} \operatorname{BER}(f) - \lambda \cdot \overline{\operatorname{BER}}(f)$$

Easy to show that

$$f^*(x) = \llbracket \eta(x) - \pi > \lambda \cdot (\bar{\eta}(x) - \bar{\pi}) \rrbracket$$

$$\eta(x) \doteq \mathbb{P}(Y = +1 \mid X = x)$$

$$\bar{\pi} \doteq \mathbb{P}(\bar{Y} = +1)$$

# CPE approach?

Alternately, let us consider the Bayes-optimal solutions

$$f^* \in \operatorname{argmin}_{f: \mathcal{X} \rightarrow \{\pm 1\}} \operatorname{BER}(f) - \lambda \cdot \overline{\operatorname{BER}}(f)$$

Easy to show that

$$f^*(x) = \llbracket \eta(x) - \pi > \lambda \cdot (\bar{\eta}(x) - \bar{\pi}) \rrbracket$$

$$\eta(x) \doteq \mathbb{P}(Y = +1 \mid X = x)$$

$$\bar{\eta}(x) \doteq \mathbb{P}(\bar{Y} = +1 \mid X = x)$$

$$\bar{\pi} \doteq \mathbb{P}(\bar{Y} = +1)$$

$$\pi \doteq \mathbb{P}(Y = +1)$$

# CPE approach?

Alternately, let us consider the Bayes-optimal solutions

$$f^* \in \operatorname{argmin}_{f: \mathcal{X} \rightarrow \{\pm 1\}} \operatorname{BER}(f) - \lambda \cdot \overline{\operatorname{BER}}(f)$$

Easy to show that

$$f^*(x) = \llbracket \eta(x) - \pi > \lambda \cdot (\bar{\eta}(x) - \bar{\pi}) \rrbracket$$

$$\eta(x) \doteq \mathbb{P}(Y = +1 \mid X = x)$$

$$\bar{\eta}(x) \doteq \mathbb{P}(\bar{Y} = +1 \mid X = x)$$

$$\bar{\pi} \doteq \mathbb{P}(\bar{Y} = +1)$$

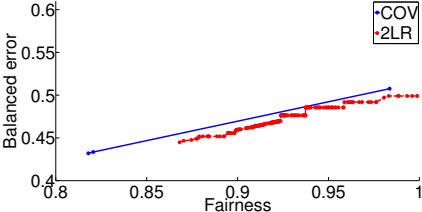
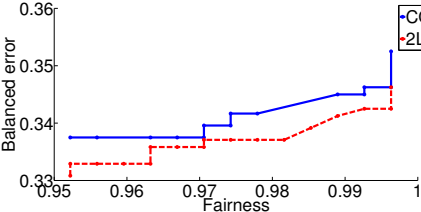
$$\pi \doteq \mathbb{P}(Y = +1)$$

**Just requires CPE on the target and sensitive features!**

- tuning of  $\lambda$  does not require re-training
- also useful to study feature learning (McNamara et al., 2017)

# Illustration of CPE approach

Competitive performance with bespoke optimisation (COV) on UCI adult and synthetic Gaussian datasets

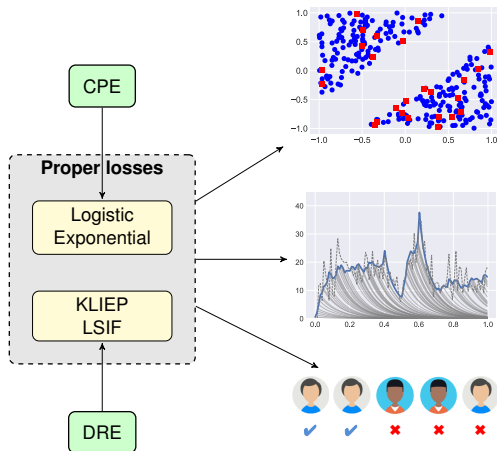


# Conclusion

# Talk summary

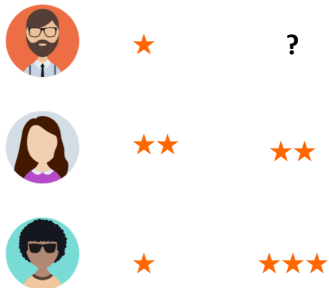
A formal link between DRE and CPE

CPE approach to three distinct learning problems

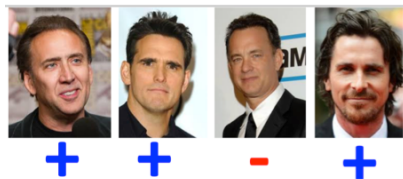




# For another day



*Recommender systems*



*Ranking*

# Collaborators



Brendan van Rooyen  
ANU



Bob Williamson  
ANU



Cheng Soon Ong  
Data61/ANU



Richard Nock  
Data61/ANU



Nagarajan Natarajan  
MSR Bangalore



Giorgio Patrini  
UvA-Bosch DELTA



Young Lee  
Data61/ANU



Lizhen Qu  
Data61/ANU

Thanks!

# Further reading

Linking losses for density ratio and class-probability estimation. Aditya Krishna Menon and Cheng Soon Ong. ICML 2016.

A scaled Bregman theorem with applications. Richard Nock, Aditya Krishna Menon and Cheng Soon Ong. NIPS 2016.

---

Learning from corrupted binary labels via class-probability estimation. Aditya Krishna Menon, Brendan van Rooyen, Cheng Soon Ong and Robert C. Williamson. ICML 2015.

Learning with symmetric label noise: the importance of being unhinged. Brendan van Rooyen, Aditya Krishna Menon and Robert C. Williamson. NIPS 2015.

Learning from binary labels with instance-dependent corruption. Aditya Krishna Menon, Brendan van Rooyen and Nagarajan Natarajan. <https://arxiv.org/abs/1605.00751>

Making deep neural networks robust to label noise: a loss correction approach. Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, Lizhen Qu. CVPR 2017.

---

Beyond the likelihood: new loss functions for (non-)linear Hawkes processes. Aditya Krishna Menon and Young Lee. In preparation.

---

The cost of fairness in binary classification. Aditya Krishna Menon and Robert C. Williamson. <https://arxiv.org/abs/1705.09055>