

COMS 4771 Lecture 13

1. Beyond prediction error
2. Beyond binary classification

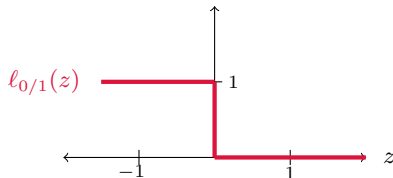
BEYOND PREDICTION ERROR

PREDICTION ERROR / ZERO-ONE LOSS

P is a distribution over $\mathcal{X} \times \{-1, +1\}$, and $(X, Y) \sim P$.

For any classifier $f: \mathcal{X} \rightarrow \{-1, +1\}$,

$$\text{err}(f) = \Pr[f(X) \neq Y] = \mathbb{E}[\ell_{0/1}(Yf(X))].$$

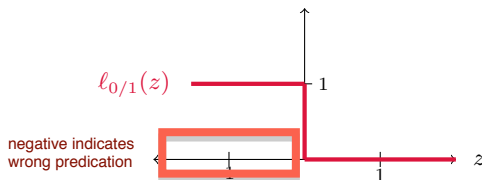


PREDICTION ERROR / ZERO-ONE LOSS

P is a distribution over $\mathcal{X} \times \{-1, +1\}$, and $(X, Y) \sim P$.

For any classifier $f: \mathcal{X} \rightarrow \{-1, +1\}$,

$$\text{err}(f) = \Pr[f(X) \neq Y] = \mathbb{E}[\ell_{0/1}(Yf(X))].$$



Also works with **real-valued predictors** $f: \mathcal{X} \rightarrow \mathbb{R}$; for example:

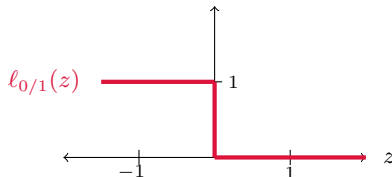
- ▶ **Linear classifiers:** $x \mapsto \langle w, x \rangle - \theta$.
- ▶ **(Binary) plug-in classifiers** $x \mapsto \widehat{\Pr}[Y = +1 \mid X = x] - 1/2$.

PREDICTION ERROR / ZERO-ONE LOSS

P is a distribution over $\mathcal{X} \times \{-1, +1\}$, and $(X, Y) \sim P$.

For any classifier $f: \mathcal{X} \rightarrow \{-1, +1\}$,

$$\text{err}(f) = \Pr[f(X) \neq Y] = \mathbb{E}[\ell_{0/1}(Yf(X))].$$



Also works with **real-valued predictors** $f: \mathcal{X} \rightarrow \mathbb{R}$; for example:

- ▶ **Linear classifiers:** $x \mapsto \langle w, x \rangle - \theta$.
- ▶ **(Binary) plug-in classifiers:** $x \mapsto \widehat{\Pr}[Y = +1 | X = x] - 1/2$.

Often useful to adjust threshold (e.g., θ and $1/2$ above).

adjust the θ to
minimize the loss function

THRESHOLDS

Uses for adjusting threshold t

Often have **different costs for different kinds of mistakes** (recall HW1):

| | $f(X) \leq t$ | $f(X) > t$ |
|----------|---------------|------------|
| $Y = -1$ | 0 | c |
| $Y = +1$ | $1 - c$ | 0 |

THRESHOLDS

Uses for adjusting threshold t

Often have **different costs for different kinds of mistakes** (recall HW1):

| | $f(X) \leq t$ | $f(X) > t$ |
|----------|---------------|------------|
| $Y = -1$ | 0 | c |
| $Y = +1$ | $1 - c$ | 0 |

Also, often interested in **different performance criteria**.

► **Precision:**

$$\Pr[Y = +1 \mid f(X) > t]$$

► **Recall** (a.k.a. **Sensitivity**, **True Positive Rate**):

$$\Pr[f(X) > t \mid Y = +1]$$

► **Specificity:**

$$\Pr[f(X) \leq t \mid Y = -1]$$

► **False Positive Rate:**

$$\Pr[f(X) > t \mid Y = -1]$$

CONDITIONAL PROBABILITY ESTIMATION

Sometimes would like real-valued predictor f to be related to the **conditional probability function** η :

$$\eta(x) := \Pr[Y = +1 \mid X = x].$$

CONDITIONAL PROBABILITY ESTIMATION

Sometimes would like real-valued predictor f to be related to the **conditional probability function** η :

$$\eta(x) := \Pr[Y = +1 \mid X = x].$$

Option #1: Can use **generative models** to construct **plug-in classifier**, but **could fail if class conditional distributions are not accurate**.

CONDITIONAL PROBABILITY ESTIMATION

motivation

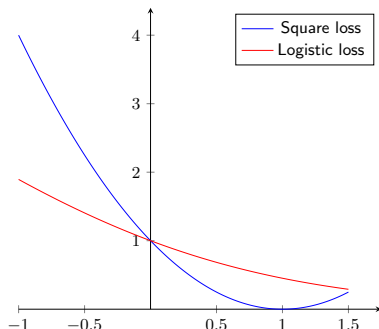
Sometimes would like real-valued predictor f to be related to the **conditional probability function** η :

$$\eta(x) := \Pr[Y = +1 \mid X = x]. \quad \text{for a random } x$$

Option #1: Can use **generative models** to construct **plug-in classifier**, but could fail if **class conditional distributions are not accurate**.

Option #2: Can use a loss function that is minimized by η (or some invertible transformation thereof).

ELICITING CONDITIONAL PROBABILITIES



Goal: loss function that is minimized by (some invertible transformation of) the conditional probability function

$$\eta(x) = \Pr[Y = +1 \mid X = x].$$

► **Square loss:** $\ell_{\text{sq}}(z) = (1 - z)^2$

$\mathbb{E}[\ell_{\text{sq}}(Y f(x)) \mid X = x]$ is minimized by $f(x) = 2\eta(x) - 1$.

► **Logistic loss:** $\ell_{\text{log}}(z) = \ln(1 + \exp(-z))$

$\mathbb{E}[\ell_{\text{log}}(Y f(x)) \mid X = x]$ is minimized by $f(x) = \ln\left(\frac{\eta(x)}{1 - \eta(x)}\right)$.

ELICITING CONDITIONAL PROBABILITIES

Using loss functions: easy with linear/affine functions whenever the loss function ℓ is a convex function (due to affine composition rule).

$$\min_{\mathbf{w} \in \mathbb{R}^d} \quad R(\mathbf{w}) + \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{x}^{(i)}, y^{(i)})$$

(Here, the regularization function R is also assumed to be convex.)

ELICITING CONDITIONAL PROBABILITIES

Using loss functions: easy with linear/affine functions whenever the loss function ℓ is a convex function (due to affine composition rule).

$$\min_{\mathbf{w} \in \mathbb{R}^d} R(\mathbf{w}) + \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{w}; \mathbf{x}^{(i)}, y^{(i)})$$

(Here, the regularization function R is also assumed to be convex.)

Good news: Both ℓ_{sq} and ℓ_{log} are convex!

ELICITING CONDITIONAL PROBABILITIES

A non-example:

- ▶ **Hinge loss:** $\ell_{\text{hl}}(z) = \max\{0, 1 - z\}$

$\mathbb{E}[\ell_{\text{hl}}(Y f(x)) | X = x]$ is minimized by $f(x) = \text{sign}(2\eta(x) - 1)$.

Can't recover $\eta(x)$ from $f(x)$.

ELICITING CONDITIONAL PROBABILITIES

A non-example:

► **Hinge loss:** $\ell_{\text{hl}}(z) = \max\{0, 1 - z\}$

$\mathbb{E}[\ell_{\text{hl}}(Y f(x)) | X = x]$ is minimized by $f(x) = \text{sign}(2\eta(x) - 1)$.

Can't recover $\eta(x)$ from $f(x)$.

Caveat: Might not be possible to represent

$$\mathbf{x} \mapsto 2\eta(\mathbf{x}) - 1 \quad \text{or} \quad \mathbf{x} \mapsto \ln\left(\frac{\eta(\mathbf{x})}{1 - \eta(\mathbf{x})}\right)$$

as (say) a linear function $\mathbf{x} \mapsto \langle \mathbf{w}, \mathbf{x} \rangle$.

ELICITING CONDITIONAL PROBABILITIES

A non-example:

- **Hinge loss:** $\ell_{\text{hl}}(z) = \max\{0, 1 - z\}$

$\mathbb{E}[\ell_{\text{hl}}(Y f(x)) | X = x]$ is minimized by $f(x) = \text{sign}(2\eta(x) - 1)$.

Can't recover $\eta(x)$ from $f(x)$.

Caveat: Might not be possible to represent

$$x \mapsto 2\eta(x) - 1 \quad \text{or} \quad x \mapsto \ln\left(\frac{\eta(x)}{1 - \eta(x)}\right)$$

as (say) a linear function $x \mapsto \langle w, x \rangle$.

Common solutions: enhance the feature space via feature expansion or kernels, or use more flexible models—discussed later.

BEYOND BINARY CLASSIFICATION

REDUCTIONS

Suppose we have a good algorithm \mathcal{A} for learning binary classifiers.

Question: Can we use \mathcal{A} to solve other machine learning problems?

REDUCTIONS

Suppose we have a good algorithm \mathcal{A} for learning binary classifiers.

Question: Can we use \mathcal{A} to solve other machine learning problems?

A **reduction** from Problem 1 to Problem 2 is specified by two procedures:

- ▶ **Instance map R :** transforms Problem 1 instance S into a Problem 2 instance S' .
- ▶ **Solution map T :** transforms solution f' for Problem 2 instance S' into a solution f for Problem 1 instance S .

REDUCTIONS

Suppose we have a good algorithm \mathcal{A} for learning binary classifiers.

Question: Can we use \mathcal{A} to solve other machine learning problems?

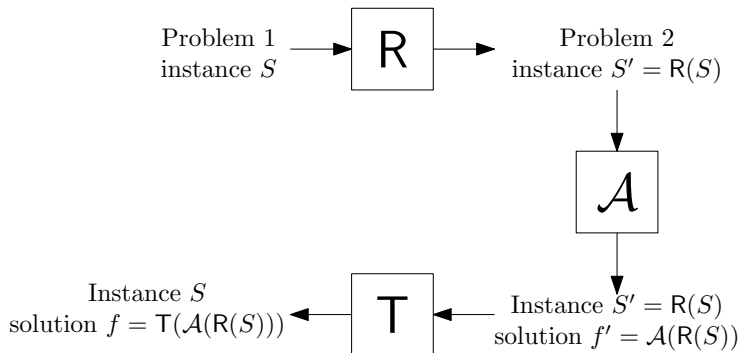
A **reduction** from Problem 1 to Problem 2 is specified by two procedures:

- ▶ **Instance map R :** transforms Problem 1 instance S into a Problem 2 instance S' .
- ▶ **Solution map T :** transforms solution f' for Problem 2 instance S' into a solution f for Problem 1 instance S .

Desired property:

If f' is a good solution for S' , then f is a good solution for S .

REDUCTIONS



In machine learning, typically have

- ▶ **Problem 1**: the problem you have to solve for a real application
- ▶ **Problem 2**: a well-studied problem in machine learning
- ▶ **Problem instance**: training data and (implicitly) a probability distribution P
- ▶ **Solution**: prediction function
- ▶ \mathcal{A} : the latest, greatest learning algorithm for Problem 2

EXAMPLES

1. **Problem:** importance-weighted classification
 - ▶ **Reduction:** rejection sampling
(Reduces problem to unweighted classification.)
2. **Problem:** multi-class classification
 - ▶ **Reduction #1:** One-Against-All
 - ▶ **Reduction #2:** Error Correcting Output Codes
(Both reduce problem to binary classification.)

IMPORTANCE-WEIGHTED CLASSIFICATION

Problem:

- ▶ **Setting:** Random triple $(X, Y, W) \sim P$ for some probability distribution P over $\mathcal{X} \times \mathcal{Y} \times \mathbb{R}_+$.

$W =$ **importance weight** for labeled example (X, Y) .

- ▶ **Goal:** Function $f: \mathcal{X} \rightarrow \mathcal{Y}$ with small **importance-weighted error**:

$$\mathbb{E}\left[W \cdot \mathbb{1}\{f(X) \neq Y\}\right].$$

IMPORTANCE-WEIGHTED CLASSIFICATION

Problem:

- ▶ **Setting:** Random triple $(X, Y, W) \sim P$ for some probability distribution P over $\mathcal{X} \times \mathcal{Y} \times \mathbb{R}_+$.

$W =$ **importance weight** for labeled example (X, Y) .

- ▶ **Goal:** Function $f: \mathcal{X} \rightarrow \mathcal{Y}$ with small **importance-weighted error**:

$$\mathbb{E} \left[W \cdot \mathbb{1}\{f(X) \neq Y\} \right].$$

Problem instance:

- ▶ Training data S : collection of triples $(x, y, w) \in \mathcal{X} \times \mathcal{Y} \times \mathbb{R}_+$, presumed to be drawn i.i.d. from P .

IMPORTANCE-WEIGHTED CLASSIFICATION

Problem:

- ▶ **Setting:** Random triple $(X, Y, W) \sim P$ for some probability distribution P over $\mathcal{X} \times \mathcal{Y} \times \mathbb{R}_+$.

$W =$ **importance weight** for labeled example (X, Y) .

- ▶ **Goal:** Function $f: \mathcal{X} \rightarrow \mathcal{Y}$ with small **importance-weighted error**:

$$\mathbb{E}\left[W \cdot \mathbb{1}\{f(X) \neq Y\}\right].$$

Problem instance:

- ▶ Training data S : collection of triples $(x, y, w) \in \mathcal{X} \times \mathcal{Y} \times \mathbb{R}_+$, presumed to be drawn i.i.d. from P .

Where it comes up:

- ▶ *Class-specific weights:* e.g., $W = 100 \Leftrightarrow Y = 0$ (and $W = 1$ otherwise).
- ▶ *Covariate-specific weights:* e.g., $W = 100 \Leftrightarrow X \in \mathcal{X}_0$ (and $W = 1$ o.w.).
- ▶ *Boosting, domain adaptation, causal inference, ...*

(Note: many learning algorithms natively handle importance weights.)

IMPORTANCE-WEIGHTED CLASSIFICATION

Problem:

- ▶ **Setting:** Random triple $(X, Y, W) \sim P$ for some probability distribution P over $\mathcal{X} \times \mathcal{Y} \times \mathbb{R}_+$.

$W =$ **importance weight** for labeled example (X, Y) .

- ▶ **Goal:** Function $f: \mathcal{X} \rightarrow \mathcal{Y}$ with small **importance-weighted error**:

$$\mathbb{E}\left[W \cdot \mathbb{1}\{f(X) \neq Y\}\right].$$

Problem instance:

- ▶ Training data S : collection of triples $(x, y, w) \in \mathcal{X} \times \mathcal{Y} \times \mathbb{R}_+$, presumed to be drawn i.i.d. from P .

Where it comes up:

- ▶ *Class-specific weights:* e.g., $W = 100 \Leftrightarrow Y = 0$ (and $W = 1$ otherwise).
- ▶ *Covariate-specific weights:* e.g., $W = 100 \Leftrightarrow X \in \mathcal{X}_0$ (and $W = 1$ o.w.).
- ▶ *Boosting, domain adaptation, causal inference, ...*

(Note: many learning algorithms natively handle importance weights.)

Would like to reduce to (unweighted) classification.

THE REJECTION SAMPLING REDUCTION

Main idea: Transform training data S so it looks like it came from a distribution P' , where

$$\mathbb{E}_{(X,Y,W) \sim P} \left[W \cdot \mathbb{1}\{f(X) \neq Y\} \right] = \mathbb{E}_{(X',Y') \sim P'} \left[\mathbb{1}\{f(X') \neq Y'\} \right].$$

THE REJECTION SAMPLING REDUCTION

Main idea: Transform training data S so it looks like it came from a distribution P' , where

$$\mathbb{E}_{(X,Y,W) \sim P} \left[W \cdot \mathbb{1}\{f(X) \neq Y\} \right] = \mathbb{E}_{(X',Y') \sim P'} \left[\mathbb{1}\{f(X') \neq Y'\} \right].$$

Instance mapping procedure

Input Training data S from $\mathcal{X} \times \mathcal{Y} \times \mathbb{R}_+$.

- 1: Initialize $S' = \emptyset$.
- 2: Let $w_{\max} := \max_{(x,y,w) \in S} w$.
- 3: **for** each $(x, y, w) \in S$ **do**
- 4: Toss a coin with heads bias $\frac{w}{w_{\max}}$.
- 5: If heads, keep example—put (x, y) into S' .
- 6: If tails, discard example.
- 7: **end for**
- 8: **return** Training data S' from $\mathcal{X} \times \mathcal{Y}$.

THE REJECTION SAMPLING REDUCTION

Main idea: Transform training data S so it looks like it came from a distribution P' , where

$$\mathbb{E}_{(X,Y,W) \sim P} \left[W \cdot \mathbb{1}\{f(X) \neq Y\} \right] = \mathbb{E}_{(X',Y') \sim P'} \left[\mathbb{1}\{f(X') \neq Y'\} \right].$$

Instance mapping procedure

Input Training data S from $\mathcal{X} \times \mathcal{Y} \times \mathbb{R}_+$.

- 1: Initialize $S' = \emptyset$.
- 2: Let $w_{\max} := \max_{(x,y,w) \in S} w$.
- 3: **for** each $(x, y, w) \in S$ **do**
- 4: Toss a coin with heads bias $\frac{w}{w_{\max}}$.
- 5: If heads, keep example—put (x, y) into S' .
- 6: If tails, discard example.
- 7: **end for**
- 8: **return** Training data S' from $\mathcal{X} \times \mathcal{Y}$.

Solution mapping procedure: identity map

THE REJECTION SAMPLING REDUCTION

Why rejection sampling works: (Assume for simplicity that $w_{\max} = 1$.)

Define random variable

$$Q := \mathbb{1}\{\text{Keep example } (X, Y)\}$$

which, after conditioning on (X, Y, W) , has mean W .

THE REJECTION SAMPLING REDUCTION

Why rejection sampling works: (Assume for simplicity that $w_{\max} = 1$.)

Define random variable

$$Q := \mathbb{1}\{\text{Keep example } (X, Y)\}$$

which, after conditioning on (X, Y, W) , has mean W .

Distribution of examples in S' is same as that of (X, Y) conditioned on $Q = 1$.

THE REJECTION SAMPLING REDUCTION

Why rejection sampling works: (Assume for simplicity that $w_{\max} = 1$.)

Define random variable

$$Q := \mathbb{1}\{\text{Keep example } (X, Y)\}$$

which, after conditioning on (X, Y, W) , has mean W .

Distribution of examples in S' is same as that of (X, Y) conditioned on $Q = 1$.

Moreover,

$$\mathbb{E}\left[Q \cdot \mathbb{1}\{f(X) \neq Y\}\right]$$

THE REJECTION SAMPLING REDUCTION

Why rejection sampling works: (Assume for simplicity that $w_{\max} = 1$.)

Define random variable

$$Q := \mathbb{1}\{\text{Keep example } (X, Y)\}$$

which, after conditioning on (X, Y, W) , has mean W .

Distribution of examples in S' is same as that of (X, Y) conditioned on $Q = 1$.

Moreover,

$$\mathbb{E}\left[Q \cdot \mathbb{1}\{f(X) \neq Y\}\right] = \mathbb{E}\left[\mathbb{E}\left[Q \cdot \mathbb{1}\{f(X) \neq Y\} \mid (X, Y, W)\right]\right]$$

THE REJECTION SAMPLING REDUCTION

Why rejection sampling works: (Assume for simplicity that $w_{\max} = 1$.)

Define random variable

$$Q := \mathbb{1}\{\text{Keep example } (X, Y)\}$$

which, after conditioning on (X, Y, W) , has mean W .

Distribution of examples in S' is same as that of (X, Y) conditioned on $Q = 1$.

Moreover,

$$\begin{aligned}\mathbb{E}\left[Q \cdot \mathbb{1}\{f(X) \neq Y\}\right] &= \mathbb{E}\left[\mathbb{E}\left[Q \cdot \mathbb{1}\{f(X) \neq Y\} \mid (X, Y, W)\right]\right] \\ &= \mathbb{E}\left[W \cdot \mathbb{1}\{f(X) \neq Y\}\right]\end{aligned}$$

THE REJECTION SAMPLING REDUCTION

Why rejection sampling works: (Assume for simplicity that $w_{\max} = 1$.)

Define random variable

$$Q := \mathbb{1}\{\text{Keep example } (X, Y)\}$$

which, after conditioning on (X, Y, W) , has mean W .

Distribution of examples in S' is same as that of (X, Y) conditioned on $Q = 1$.

Moreover,

$$\begin{aligned}\mathbb{E}\left[Q \cdot \mathbb{1}\{f(X) \neq Y\}\right] &= \mathbb{E}\left[\mathbb{E}\left[Q \cdot \mathbb{1}\{f(X) \neq Y\} \mid (X, Y, W)\right]\right] \\ &= \mathbb{E}\left[W \cdot \mathbb{1}\{f(X) \neq Y\}\right]\end{aligned}$$

Conclusion:

Prediction error w.r.t. $P' =$ importance-weighted error w.r.t. P .

MULTI-CLASS CLASSIFICATION

Problem:

- ▶ **Setting:** Random pair $(X, Y) \sim P$ for some probability distribution P over $\mathcal{X} \times \{1, 2, \dots, K\}$.
- ▶ **Goal:** Function $f: \mathcal{X} \rightarrow \mathcal{Y}$ with small prediction error $\Pr(f(X) \neq Y)$.

MULTI-CLASS CLASSIFICATION

Problem:

- ▶ **Setting:** Random pair $(X, Y) \sim P$ for some probability distribution P over $\mathcal{X} \times \{1, 2, \dots, K\}$.
- ▶ **Goal:** Function $f: \mathcal{X} \rightarrow \mathcal{Y}$ with small prediction error $\Pr(f(X) \neq Y)$.

Problem instance:

- ▶ Training data S : collection of pairs $(x, y) \in \mathcal{X} \times \{1, 2, \dots, K\}$, presumed to be drawn i.i.d. from P .

MULTI-CLASS CLASSIFICATION

Problem:

- ▶ **Setting:** Random pair $(X, Y) \sim P$ for some probability distribution P over $\mathcal{X} \times \{1, 2, \dots, K\}$.
- ▶ **Goal:** Function $f: \mathcal{X} \rightarrow \mathcal{Y}$ with small prediction error $\Pr(f(X) \neq Y)$.

Problem instance:

- ▶ Training data S : collection of pairs $(x, y) \in \mathcal{X} \times \{1, 2, \dots, K\}$, presumed to be drawn i.i.d. from P .

Would like to reduce to binary classification.

ONE-AGAINST-ALL REDUCTION

Main idea: Create K binary classification problems

given $x \in \mathcal{X}$, predict whether or not $y = i$.

Create K examples from each $(x, y) \in S$:

$$(x, y) \longrightarrow \left\{ \begin{array}{lll} (x, \mathbb{1}\{y = 1\}) & \longrightarrow & S'_1 \\ (x, \mathbb{1}\{y = 2\}) & \longrightarrow & S'_2 \\ \vdots & \vdots & \vdots \\ (x, \mathbb{1}\{y = K\}) & \longrightarrow & S'_K \end{array} \right.$$

ONE-AGAINST-ALL REDUCTION

Instance mapping procedure

Input Training data S from $\mathcal{X} \times \{1, 2, \dots, K\}$.

- 1: Initialize empty sets S'_1, S'_2, \dots, S'_K .
- 2: **for** each $(x, y) \in S$ **do**
- 3: **for** each $i = 1, 2, \dots, K$ **do**
- 4: Put $(x, \mathbb{1}\{y = i\}) \in \mathcal{X} \times \{0, 1\}$ into S'_i .
- 5: **end for**
- 6: **end for**
- 7: **return** Training data sets S'_1, S'_2, \dots, S'_K from $\mathcal{X} \times \{0, 1\}$.

ONE-AGAINST-ALL REDUCTION

Instance mapping procedure

Input Training data S from $\mathcal{X} \times \{1, 2, \dots, K\}$.

- 1: Initialize empty sets S'_1, S'_2, \dots, S'_K .
- 2: **for** each $(x, y) \in S$ **do**
- 3: **for** each $i = 1, 2, \dots, K$ **do**
- 4: Put $(x, \mathbb{1}\{y = i\}) \in \mathcal{X} \times \{0, 1\}$ into S'_i .
- 5: **end for**
- 6: **end for**
- 7: **return** Training data sets S'_1, S'_2, \dots, S'_K from $\mathcal{X} \times \{0, 1\}$.

Solution mapping procedure

Input K binary predictors $f'_1, f'_2, \dots, f'_K: \mathcal{X} \rightarrow \{0, 1\}$.

return Function $f: \mathcal{X} \rightarrow \{1, 2, \dots, K\}$ where

$$f(x) = \arg \max_{i \in \{1, 2, \dots, K\}} f'_i(x) \quad (\text{breaking ties arbitrarily}).$$

PROBLEM WITH OAA

OAA multi-class predictor:

$$f(x) = \arg \max_{i \in \{1, 2, \dots, K\}} f'_i(x).$$

PROBLEM WITH OAA

OAA multi-class predictor:

$$f(x) = \arg \max_{i \in \{1, 2, \dots, K\}} f'_i(x).$$

Only get correct classification on (x, y) if $f'_y(x) = 1$ and $f'_i(x) = 0$ for all $i \neq y$.
(Could err if any of the f'_i errs!)

PROBLEM WITH OAA

OAA multi-class predictor:

$$f(x) = \arg \max_{i \in \{1, 2, \dots, K\}} f'_i(x).$$

Only get correct classification on (x, y) if $f'_y(x) = 1$ and $f'_i(x) = 0$ for all $i \neq y$.
(Could err if any of the f'_i errs!)

Common solution: use conditional probability estimation

$$f'_i(x) = \text{estimate of } \Pr[Y = i \mid X = x].$$

ERROR CORRECTING OUTPUT CODES

Main idea: Create m binary classification problems

given $x \in \mathcal{X}$, predict whether or not $y \in L_i$

for m pre-specified subsets $L_i \subset \{1, 2, \dots, K\}$.

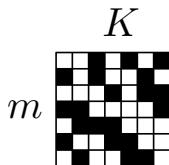
ERROR CORRECTING OUTPUT CODES

Main idea: Create m binary classification problems

given $x \in \mathcal{X}$, predict whether or not $y \in L_i$

for m pre-specified subsets $L_i \subset \{1, 2, \dots, K\}$.

Subsets specified using *error correcting codes* ($m \times K$ binary matrix):



Let i -th column $\mathbf{c}_i \in \{0, 1\}^m$ be the *code word* for class $i \in \{1, 2, \dots, K\}$.

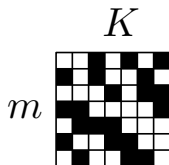
ERROR CORRECTING OUTPUT CODES

Main idea: Create m binary classification problems

given $x \in \mathcal{X}$, predict whether or not $y \in L_i$

for m pre-specified subsets $L_i \subset \{1, 2, \dots, K\}$.

Subsets specified using *error correcting codes* ($m \times K$ binary matrix):



Let i -th column $\mathbf{c}_i \in \{0, 1\}^m$ be the *code word* for class $i \in \{1, 2, \dots, K\}$.

Use m binary classifiers f'_1, f'_2, \dots, f'_m to *predict* entire code word.

ERROR CORRECTING OUTPUT CODES

Instance mapping procedure: similar to OAA.

ERROR CORRECTING OUTPUT CODES

Instance mapping procedure: similar to OAA.

Solution mapping procedure:

Input m binary predictors $f'_1, f'_2, \dots, f'_m: \mathcal{X} \rightarrow \{0, 1\}$.

return Function $f: \mathcal{X} \rightarrow \{1, 2, \dots, K\}$ where

$$f(x) = \arg \min_{i \in \{1, 2, \dots, K\}} \left\| \begin{bmatrix} f'_1(x) \\ f'_m(x) \\ \vdots \\ f'_m(x) \end{bmatrix} - \mathbf{c}_i \right\|.$$

ERROR CORRECTING OUTPUT CODES

Instance mapping procedure: similar to OAA.

Solution mapping procedure:

Input m binary predictors $f'_1, f'_2, \dots, f'_m: \mathcal{X} \rightarrow \{0, 1\}$.

return Function $f: \mathcal{X} \rightarrow \{1, 2, \dots, K\}$ where

$$f(x) = \arg \min_{i \in \{1, 2, \dots, K\}} \left\| \begin{bmatrix} f'_1(x) \\ f'_m(x) \\ \vdots \\ f'_m(x) \end{bmatrix} - \mathbf{c}_i \right\|.$$

Error correction: can still get correct multi-class prediction even if several binary classifiers err.

ERROR CORRECTING OUTPUT CODES

Instance mapping procedure: similar to OAA.

Solution mapping procedure:

Input m binary predictors $f'_1, f'_2, \dots, f'_m: \mathcal{X} \rightarrow \{0, 1\}$.

return Function $f: \mathcal{X} \rightarrow \{1, 2, \dots, K\}$ where

$$f(x) = \arg \min_{i \in \{1, 2, \dots, K\}} \left\| \begin{bmatrix} f'_1(x) \\ f'_m(x) \\ \vdots \\ f'_m(x) \end{bmatrix} - \mathbf{c}_i \right\|.$$

Error correction: can still get correct multi-class prediction even if several binary classifiers err.

Still also useful to use conditional probability estimation

$$f'_i(x) = \text{estimate of } \Pr[Y \in L_i \mid X = x]$$

("Probabilistic ECOC").

ERROR CORRECTING OUTPUT CODES

Instance mapping procedure: similar to OAA.

Solution mapping procedure:

Input m binary predictors $f'_1, f'_2, \dots, f'_m: \mathcal{X} \rightarrow \{0, 1\}$.

return Function $f: \mathcal{X} \rightarrow \{1, 2, \dots, K\}$ where

$$f(x) = \arg \min_{i \in \{1, 2, \dots, K\}} \left\| \begin{bmatrix} f'_1(x) \\ f'_m(x) \\ \vdots \\ f'_m(x) \end{bmatrix} - \mathbf{c}_i \right\|.$$

Error correction: can still get correct multi-class prediction even if several binary classifiers err.

Still also useful to use conditional probability estimation

$$f'_i(x) = \text{estimate of } \Pr[Y \in L_i \mid X = x]$$

(“Probabilistic ECOC”).

Caveat: binary prediction problems could be challenging / unnatural (e.g., predict if handwritten digit is an even digit or not).

COMPARING OAA AND ECOC

Empirical comparison

- ▶ Eight multi-class problems (from the UCI repository).
- ▶ $\mathcal{A} = \text{classregtree}$ from the MATLAB statistics toolbox, estimate conditional probabilities using square loss.
- ▶ ECOC based on Hadamard matrix (similar to Fourier transform).

| Data set | Number of classes | One-against-all | ECOC |
|-----------|-------------------|-----------------|---------------|
| ecoli | 8 | 0.0985 | 0.0517 |
| glass | 6 | 0.3874 | 0.3462 |
| pendigits | 10 | 0.0985 | 0.0517 |
| satimage | 6 | 0.1679 | 0.1376 |
| soybean | 19 | 0.6580 | 0.5993 |
| splice | 3 | 0.0642 | 0.0699 |
| vowel | 11 | 0.6356 | 0.5780 |
| yeast | 10 | 0.4893 | 0.4479 |

- ▶ **Reductions:** reuse existing technology to solve new problems.
 - ▶ Multi-class (OAA, ECOC, tournaments, ...)
 - ▶ Multi-label prediction
 - ▶ Ranking
 - ▶ Sequence prediction
 - ▶ ...
- ▶ **Lots of different problems and objectives** beyond binary classification and prediction error—can be application-/domain-specific.
- ▶ **Very useful tools:**
 - ▶ Importance-weighted classification (using reductions)
 - ▶ Conditional probability estimation (using certain loss functions)