

# Supervised Machine Learning Review

## Outline

by Paul Hand  
Northeastern University

Regression + Classification Problems

Statistical Framework for ML

Justification for square loss & cross entropy loss

Bias Variance Trade off, model selection, an unexpected twist

## Common Problems in Supervised ML

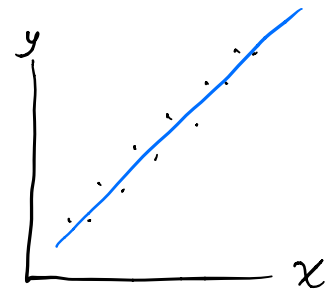
**Regression**: predict a continuous value

Let  $f: \mathbb{R}^d \rightarrow \mathbb{R}$

$$y = f(x) + \text{noise}$$

Given:  $\{(x_i, y_i)\}_{i=1 \dots n}$

Find:  $f$



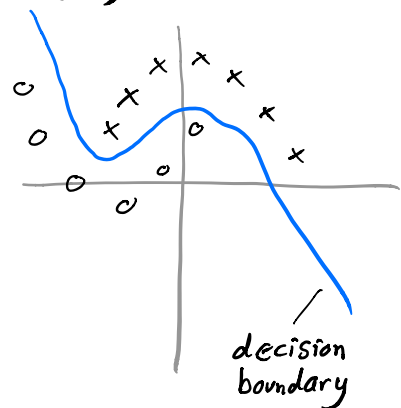
**Classification**: predict membership in a category

Let  $f: \mathbb{R}^d \rightarrow \begin{Bmatrix} \text{cat } 1 \\ \vdots \\ \text{cat } m \end{Bmatrix}$

$$y = f(x) + \text{noise}$$

Given:  $\{(x_i, y_i)\}_{i=1 \dots n}$

Find:  $f$



## Terminology

$x$  - input variables, predictors, independent vars, features

$y$  - response, dependent variable, output variable

$f$  - model, predictor, hypothesis

## Statistical Framework for ML (supervised)

Assume:

- $(X, y)$  are sampled from a joint probability distribution
- Training data  $D = \{(x_i, y_i)\}_{i=1 \dots n}$  are iid samples
- Test data are also iid samples

Can estimate the model/predictor by maximum likelihood estimation

Results (usually) in an optimization problem

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{H}} \sum_{i=1}^n \ell(f(x_i), y_i)$$

where

$\ell$  - loss function eg  $\ell(\hat{y}, y) = |\hat{y} - y|^2$

$\mathcal{H}$  - hypothesis class eg degree  $d$  polynomial

Q's: What loss do you choose and why?

What hypotheses should you search over?

## Linear Regression and Square Loss

Let  $a \in \mathbb{R}^d$ ,  $x \in \mathbb{R}^d$

Model:  $y_i = x_i^t a + \varepsilon_i$  w/  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$

Data:  $\mathcal{D} = \{(x_i, y_i)\}_{i=1 \dots n}$

Estimate  $a$  by maximum likelihood

pdf of  $\varepsilon_i$  is  $\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{z^2}{2\sigma^2}}$  over  $z \in \mathbb{R}$

likelihood of data (using  $\varepsilon_i = y_i - x_i^t a$ )

$$L(a) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - x_i^t a)^2}{2\sigma^2}}$$

$$\log L(a) = -\sum_{i=1}^n \frac{(y_i - x_i^t a)^2}{2\sigma^2} + \text{terms constant in } a$$

maximizing data likelihood  $\Leftrightarrow$  minimizing square loss

$$\max_a L(a) \Leftrightarrow \min_a \sum_{i=1}^n \underbrace{(x_i^t a - y_i)^2}_{\text{Square loss } \ell(\hat{y}, y) = |\hat{y} - y|^2}$$

# Logistic Regression and Cross Entropy Loss

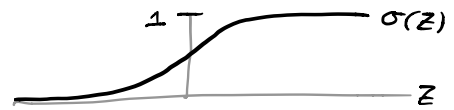
Model:

Let  $a \in \mathbb{R}^d$

$$P(y=1|x) = \sigma(x^t a)$$

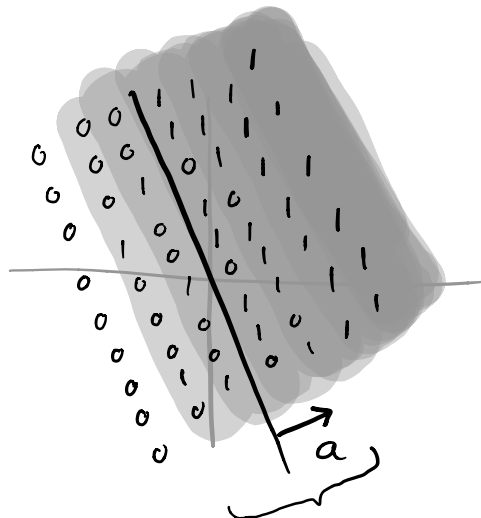
$$P(y=0|x) = 1 - \sigma(x^t a)$$

$$\text{w/ } \sigma(z) = \frac{e^z}{e^z + 1} = \frac{1}{1 + e^{-z}}$$



Data:  $\{(x_i, y_i)\}$

Visually:



width of region of uncertainty  $\approx \frac{1}{\|a\|_2}$

Estimate  $a$  by maximum likelihood

$$L(a) = \prod_{i=1}^n P(y_i=0|x_i)^{1-y_i} P(y_i=1|x_i)^{y_i}$$

$$\log L(a) = \sum_{i=1}^n (1-y_i) \log P(y_i=0|x_i) + y_i \log P(y_i=1|x_i)$$

Cross entropy loss

$$\mathcal{L}_{CE}(P, q) = - \sum_{z \in \mathcal{Z}} P(z) \log q(z) = - \mathbb{E}_P(\log q)$$

discrete  
r.v.s over  $\mathcal{Z}$

Maximizing data likelihood  $\Leftrightarrow$  minimizing cross entropy loss

$$\max_a L(a) \Leftrightarrow \min_a \underbrace{- \sum_{i=1}^n (y_i \log(\sigma(x_i^T a)) + (1-y_i) \log(1-\sigma(x_i^T a)))}_{\mathcal{L}_{CE}\left(\begin{pmatrix} y_i \\ 1-y_i \end{pmatrix}, \begin{pmatrix} \sigma(x_i^T a) \\ 1-\sigma(x_i^T a) \end{pmatrix}\right)}$$

Note:

Cross entropy loss penalizes data points of a observed category to which the model assigns a very low probability.

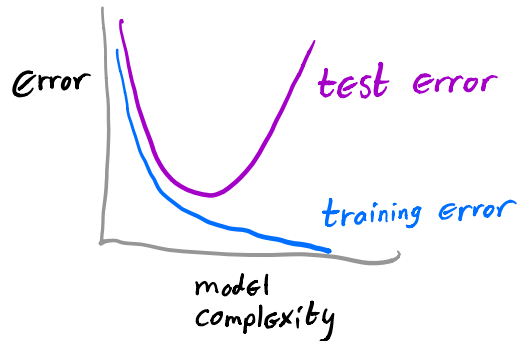
Question to ponder:

Is minimizing Cross Entropy loss all that different from minimizing a square loss in the case of regression?

# Bias-Variance Tradeoff

What class of hypotheses should you search over?

Standard Statistical ML story:



higher complexity models have lower bias but higher variance

If complexity is too high, it overfits data, variance term dominates test error

after a certain threshold, "larger models are worse"

## Bias-Variance Decomposition

Consider regression model

$$y = f(x) + \varepsilon \quad \text{w/} \quad \mathbb{E}[\varepsilon | x] = 0$$

Let  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  be iid samples

Estimate  $f$  by an algorithm producing  $\hat{f}_{\mathcal{D}}$

Evaluate  $\hat{f}_{\mathcal{D}}$  by expected loss on a new sample

$$R(\hat{f}_{\mathcal{D}}) = \mathbb{E}_{x,y} (\hat{f}_{\mathcal{D}}(x) - y)^2$$

risk                      test sample                      square loss

Performance will vary based on  $\mathcal{D}$ . Take expectation over  $\mathcal{D}$ .

$$\mathbb{E}_{\mathcal{D}} R(\hat{f}_{\mathcal{D}}) = \mathbb{E}_{x,y,\mathcal{D}} (\hat{f}_{\mathcal{D}}(x) - y)^2$$

We will decompose into 3 effects: bias, variance, irreducible error

$$\begin{aligned} \mathbb{E}_D R(\hat{f}_D) &= \mathbb{E}_{x,y,D} \left[ (\hat{f}_D(x) - f(x) - \varepsilon)^2 \right] \\ &= \mathbb{E}_{x,y,D} (\hat{f}_D(x) - f(x))^2 - 2 \mathbb{E}[(\hat{f}_D(x) - f(x))\varepsilon] + \mathbb{E}[\varepsilon^2] \\ &= \mathbb{E}_{x,y,D} (\hat{f}_D(x) - f(x))^2 + \text{Var}(\varepsilon) \end{aligned}$$

Evaluating the first term, Conditioning on  $x$ ,

$$\begin{aligned} \mathbb{E}_D (\hat{f}_D(x) - f(x))^2 &= \mathbb{E}_D \left[ \left( (\hat{f}_D(x) - \mathbb{E}_D \hat{f}_D(x)) + (\mathbb{E}_D \hat{f}_D(x) - f(x)) \right)^2 \right] \\ &= \mathbb{E}_D (\hat{f}_D(x) - \mathbb{E}_D \hat{f}_D(x))^2 + 2 \mathbb{E}_D (\hat{f}_D(x) - \mathbb{E}_D \hat{f}_D(x)) (\mathbb{E}_D \hat{f}_D(x) - f(x)) + \mathbb{E}_D (\mathbb{E}_D \hat{f}_D(x) - f(x))^2 \\ &= \underbrace{\mathbb{E}_D (\hat{f}_D(x) - \mathbb{E}_D \hat{f}_D(x))^2}_{\text{Variance of } \hat{f}_D(x)} + \underbrace{(\mathbb{E}_D (\hat{f}_D(x) - f(x)))^2}_{\text{Squared bias}} \end{aligned}$$

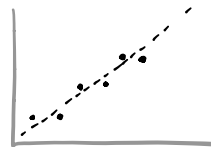
$\circ$  is expectation in  $D$       does not depend on  $D$

So,

$$\mathbb{E}_D R(\hat{f}) = \underbrace{\mathbb{E}_x (f(x) - \mathbb{E}_D \hat{f}_D(x))^2}_{\text{expected squared bias of estimate}} + \underbrace{\mathbb{E}_x \text{Var}_D \hat{f}_D(x)}_{\text{expected variance of estimate}} + \underbrace{\text{Var}(\varepsilon)}_{\text{irreducible error}}$$

Illustration of bias variance tradeoff

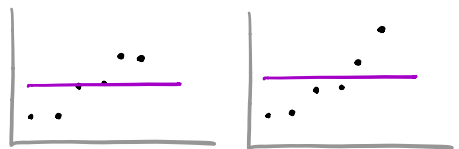
Suppose  $y = x + \varepsilon$



Low complexity model:  $y = c$

$\mathbb{E}_x (f(x) - \mathbb{E}_D \hat{f}_D)^2$  is high

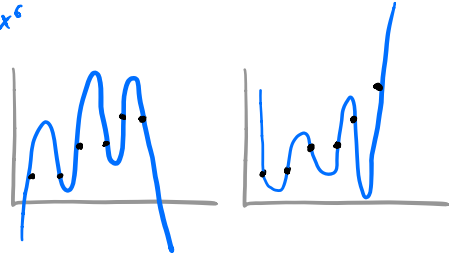
$\mathbb{E}_x \text{Var}_D \hat{f}_D(x)$  is low



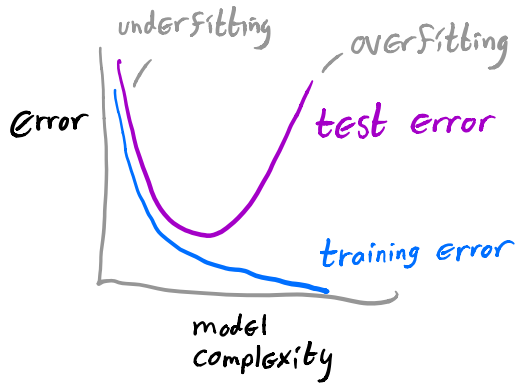
High complexity model  $\circ y = c_0 + c_1x + c_2x^2 + \dots + c_nx^n$

$\mathbb{E}_x (f(x) - \mathbb{E}_D \hat{f}_D)^2$  is low

$\mathbb{E}_x \text{Var}_D \hat{f}_D(x)$  is high



Standard Statistical ML story  $\circ$

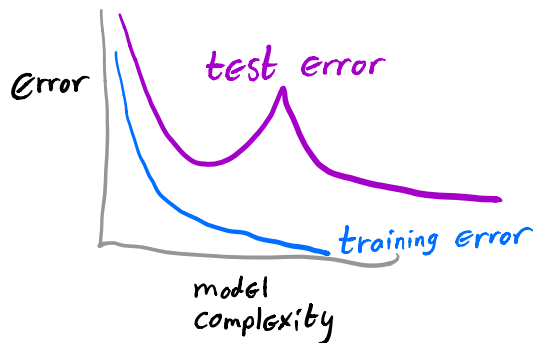


higher complexity models have lower bias but higher variance

If complexity is too high, it overfits data, variance term dominates test error

after a certain threshold, "larger models are worse"

Modern Story based on Neural Nets  $\circ$



Test error can decrease as model complexity continues increasing.

And it can be lower than in underparameterized regime

Phenomenon: double descent

underparameterized regime      overparameterized regime

"larger models are better"