# Supervised Machine Learning Review

## Outline

by Paul Hand
Northeastern University

Regression + Classification Problems

Statistical Framework for ML

Justification for Square loss & Cross entropy loss

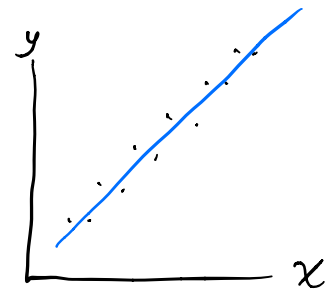Bias Variance Trade off, model selection, an unexpected twist

## Common Problems in Supervised ML

**Regression :** predict a continuous value

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

$$y = f(x) + noise$$

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

Find : $f$

**Classification:** predict membership in a category

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

$$y = f(x) + noise$$

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

Find : $f$

decision boundary

**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 $D = \{(x_i, y_i)\}_{i=1 \cdots n}$ are **iid samples**
- Test data are also **iid samples**

Can estimate the model/predictor by **maximum likelihood estimation**

Results (usually) in an optimization problem

$$\hat{f} = \underset{f \in \mathcal{H}}{\text{argmin}} \sum_{i=1}^{n} \ell(f(x_i), y_i) \qquad \text{"empirical risk minimization"}$$

$\text{separable}$

where

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

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

## What is MLE?

Estimate parameters of a model by maximizing likelihood of the observed data

## What is MLE in contrast to?

MAP - maximum a posteriori estimation - parameters have some prior distribution, data is collected, that changes the posterior distribution via Bayes Rule. Seek mode of that posterior

I just choose to minimize square loss for a binary classification problem

## Is ERM guaranteed to give you a "good" predictor?

Perhaps you get a local minimum instead of the global minimum

No, You may be doing well on training data but not on test data - overfitting

## What property is desired in the learned predictor?

Good performance on test data (future i.i.d. Samples of the distribution)

Want: Minimize the expected loss under the test distribution

## What is risk?

Risk is expected loss

## What makes $\hat{f} = \underset{f \in \mathcal{H}}{\arg\min} \sum_{i=1}^{n} \ell(f(x_i), y_i)$ empirical risk minimization?

$$\approx \underset{x,y \sim D}{\mathbb{E}} \ell(f(x), y)$$

this is empirical because it uses empirical data to estimate the expectation of loss over the training distribution

Is risk minimization biased (in a cultural sense) when applied to real problems where $\{x_i\}$ correspond to people?

Biased toward the training data - If a group is underrepresented in training data, then performance on that group may be worse

The data itself could have historical biases baked in

Just because a group has a larger fraction of the data might not mean that we want improvements in performance on that group to balance decreases in performance of other smaller groups

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/ $\varepsilon_i \sim N(0, \sigma^2)$

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

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 $\varepsilon_i = y_i - x_i^t a$)

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

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

maximizing data likelihood $\Longleftrightarrow$ minimizing square loss

$$\max_a L(a) \quad \Longleftrightarrow \quad \min_a \sum_{i=1}^{n} \underbrace{(x_i^t a - y_i)^2}_{}$$

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

# Logistic Regression and Cross Entropy Loss
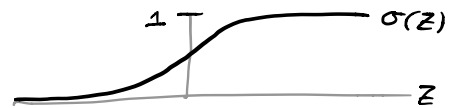
Model:

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

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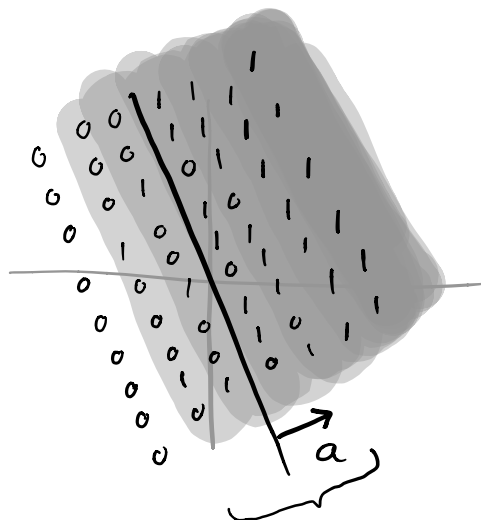$$P(y=0|x) = 1 - \sigma(x^t a)$$

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



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

$x^t a$ is a _logit_

Visually:



width of region of uncertainty $\simeq \dfrac{1}{\|a\|_2}$

Estimate $a$ by maximum likelihood

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

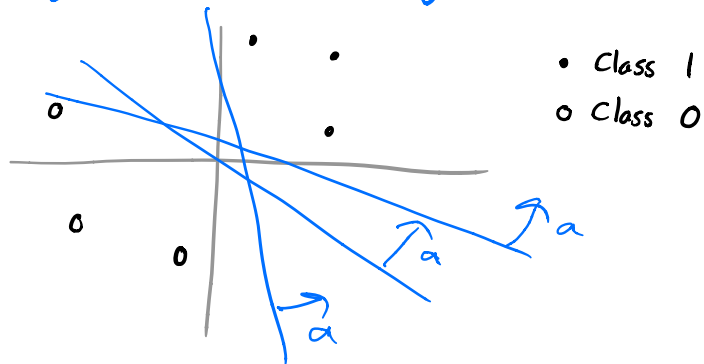$$\log L(a) = \sum_{i=1}^{n} (1-y_i) \log P(y_i=0|x_i) + y_i \log P(y_i=1|x_i)$$

Cross entropy loss

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

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

Maximizing data likelihood $\Longleftrightarrow$ minimizing cross entropy loss

$$\max_{a} L(a) \Longleftrightarrow \min_{a} -\sum_{i=1}^{n} \left( y_i \log(\sigma(x_i^t a)) + (1-y_i) \log(1-\sigma(x_i^t a)) \right)$$

$$\ell_{CE}\left( \begin{pmatrix} y_i \\ 1-y_i \end{pmatrix}, \begin{pmatrix} \sigma(x_i^t a) \\ 1-\sigma(x_i^t a) \end{pmatrix} \right)$$
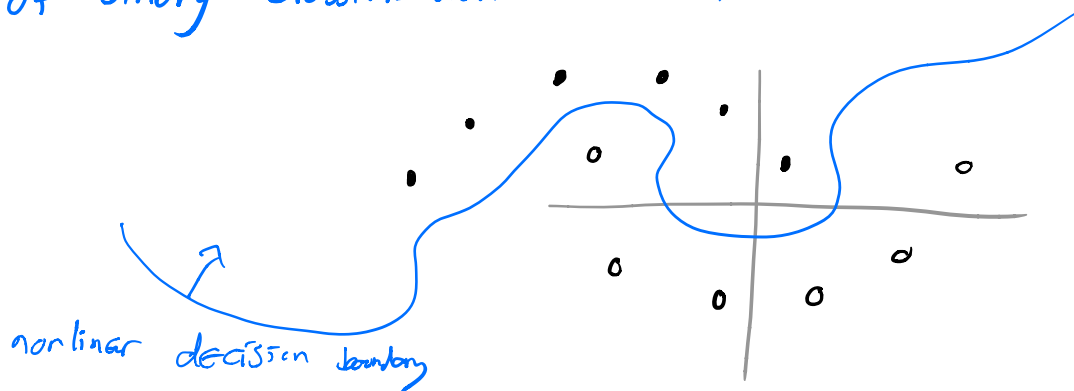
What happens if you do logistic regression on the following data by solving the above optimization problem?



• Class 1
o Class 0

What magnitude of a will result from solving this problem —- infinity - because that will increase the likelihood of the day

Cross - entropy is an asymmetric measure of the distance between two distributions.

Logistic Regression is like a simple version of binary classification w/ neural nets

nonlinear decision boundary

Note:
Cross Entropy loss penalizes data points of a observed category to which the model assigns a very low probability.
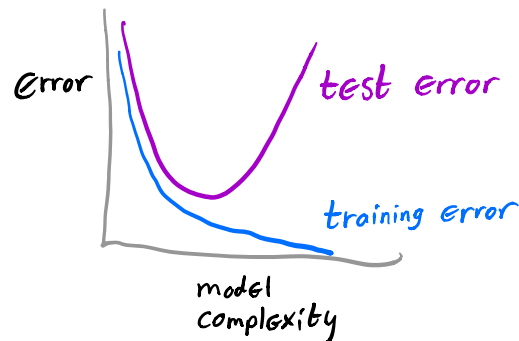
Question to ponder:
Is minimizing Cross Entropy loss all that different from minimizing a square loss in the case of logistic regression?

# Bias - Variance Tradeoff

What class of hypotheses should you search over?

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"

## Why is training error monotonically decreasing?

The search space of larger complexity models is larger

## Why is test error initially decreasing?

If its too low, it underfits the data (can not represent the "true model")

If you have $10^3$ data samples, how complex of a data model would you consider?

< 10^3.  ..... so choose something like like 30 or 100

Why does understanding this tradeoff matter?

Help select the right level of complexity

Say to look for evidence of overfitting

# Bias - Variance Decomposition

Consider regression model

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

Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1\cdots n}$ 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

$$\underset{risk}{R(\hat{f}_{\mathcal{D}})} = \underset{\substack{test \\ sample}}{\mathbb{E}_{x,y}} \; (\hat{f}_{\mathcal{D}}(x) - y)^2 \quad \underset{square\ loss}{}$$

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

$$\mathbb{E}_{\mathcal{D}} \; R(\hat{f}_{\mathcal{D}}) = \mathbb{E}_{x,y,\mathcal{D}} \left[ (\hat{f}_{\mathcal{D}}(x) - f(x) - \varepsilon)^2 \right]$$

$$= \mathbb{E}_{x,y,\mathcal{D}} (\hat{f}_{\mathcal{D}}(x) - f(x))^2 - 2 \, \mathbb{E}\left[ (\hat{f}_{\mathcal{D}}(x) - f(x))\varepsilon \right] + \mathbb{E}[\varepsilon^2]$$

$$\underbrace{\qquad\qquad}_{Var(\varepsilon)}$$

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

Evaluating the first term, Conditioning on $x$,

$$\mathbb{E}_{\mathcal{D}} (\hat{f}_{\mathcal{D}}(x) - f(x))^2 = \mathbb{E}_{\mathcal{D}} \left[ \left( (\hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}\hat{f}_{\mathcal{D}}(x)) + (\mathbb{E}_{\mathcal{D}}\hat{f}_{\mathcal{D}}(x) - f(x)) \right)^2 \right]$$

$$= \mathbb{E}_{\mathcal{D}} \left( \hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}\hat{f}_{\mathcal{D}}(x) \right)^2 + 2 \, \mathbb{E}_{\mathcal{D}} \left( \hat{f}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}\hat{f}_{\mathcal{D}}(x) \right) \left( \mathbb{E}_{\mathcal{D}}\hat{f}_{\mathcal{D}}(x) - f(x) \right) + \mathbb{E}_{\mathcal{D}} \left( \mathbb{E}_{\mathcal{D}}\hat{f}_{\mathcal{D}}(x) - f(x) \right)^2$$

$$\underbrace{\qquad\qquad}_{\substack{0\ in\ expectation \\ in\ \mathcal{D}}} \qquad \underbrace{\qquad\qquad}_{\substack{does\ not\ depend \\ on\ \mathcal{D}}}$$

$$= \underbrace{\mathbb{E}_D \left( \hat{f}_D(x) - \mathbb{E}_D \hat{f}_D(x) \right)^2}_{\text{Variance of } \hat{f}_D(x)} + \underbrace{\left( \mathbb{E}_D \left( \hat{f}_D(x) - f(x) \right) \right)^2}_{\text{squared bias}}$$

So,

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

Illustration of bias variance tradeoff

Suppose $\quad y = x + \varepsilon$



Low complexity model: $y = c$

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

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



High complexity model: $y = c_0 + c_1 x + c_2 x^2 + \cdots c_6 x^6$

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

$\quad \mathbb{E}_x \operatorname{Var}_D \hat{f}_D(x)$ is high
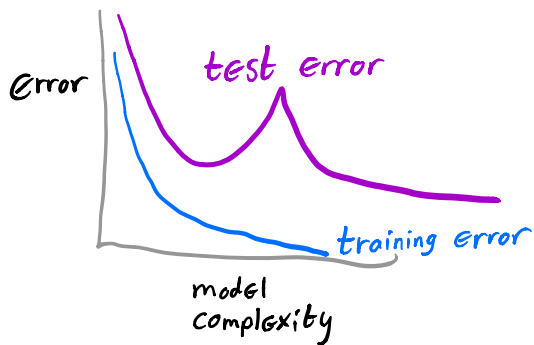
## Standard Statistical ML Story:



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"

## Modern Story based on Neural Nets:



Test error can decrease as
model complexity continues increasing.

And it can be lower than in
underparameterized regime

Phenomenon: double descent

"larger models are better"

Q: Are larger models better
b/c we have so much data that
it captures the entire problem domain

*and is actually overfitting?*

**If you have $10^3$ data samples, how complex of a data model would you consider?**

Choose a neural network with 10000 or 100000 parameters

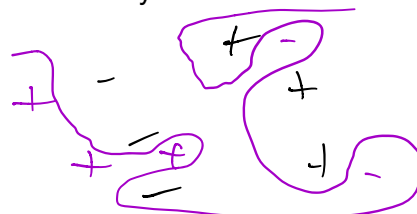**Why is being critically parameterized bad for generalization?**

Critically parameterized:  # parameters = # data points

How many values of parameters would fit data exactly? 1.   Neural net must contort itself to fit the exact data.  No expectation for generalization.

**In the overparameterized regime, do all models with 0 training error generalize well?**

there is an infinity of model parameters that fit data exactly.  Gradient descent will find one of them.  Would all solutions generalize well?

There are solutions that don't generalize well.
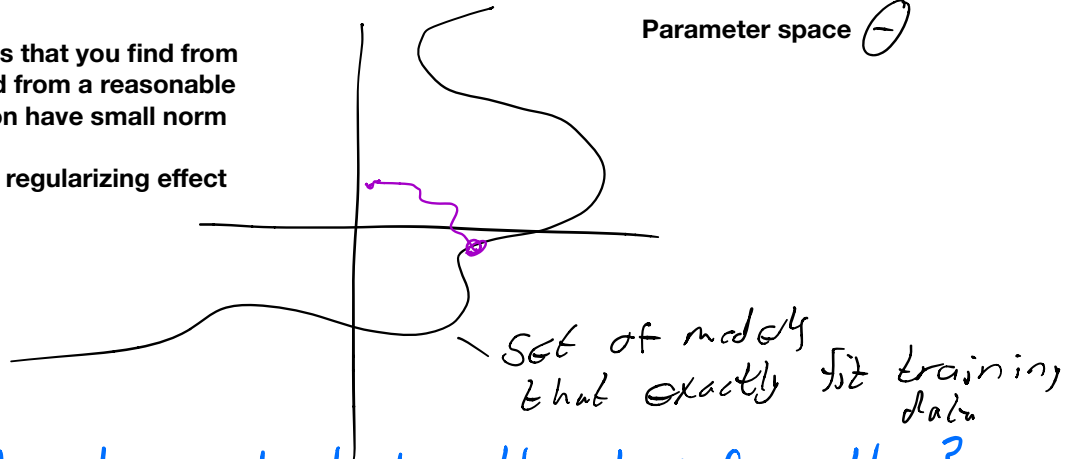Build them by adding poison training data

# How is good generalization possible in the over parameterized regime?

**Parameters that you find from running Gd from a reasonable initialization have small norm**

**That has a regularizing effect**

Parameter space

Set of models that exactly fit training data

# Why does understanding this tradeoff matter?

Expect near perfect fitting of your training data