

Tuesday, Feb 17

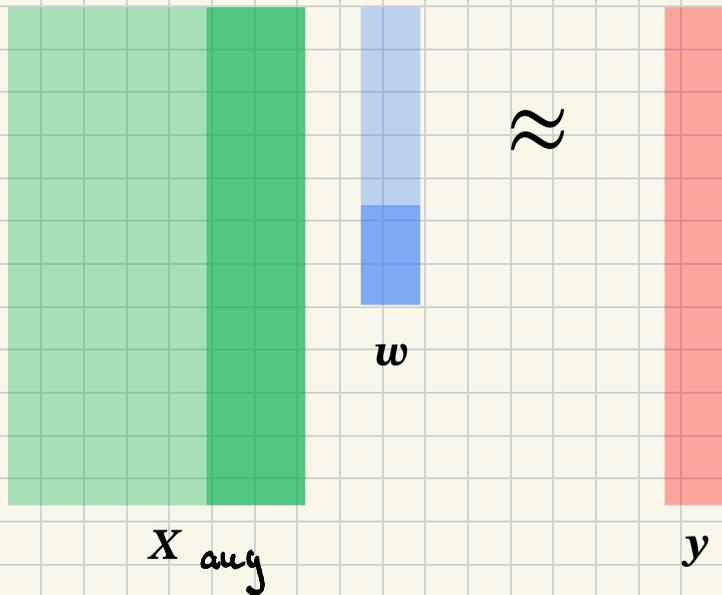
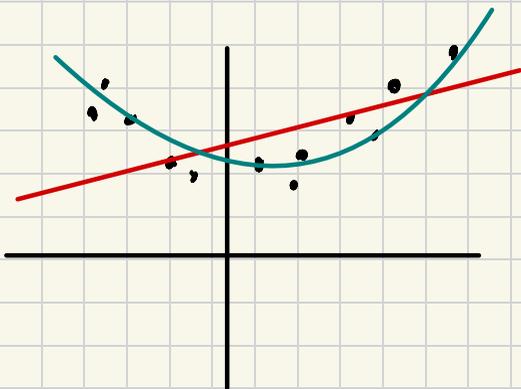
- Exam on 3/5

Plan

- Polynomial regression
- Logistic regression

Polynomial Regression

Can we do better than a linear fit?



$$f(x) = w_0 + w_1 x_1 + \dots + w_d x_d \quad (\text{linear})$$

$$f(x) = w_0 + w_1 x_1 + w_2 x_1 x_2 + w_3 x_2^2 \quad (\text{polynomial})$$

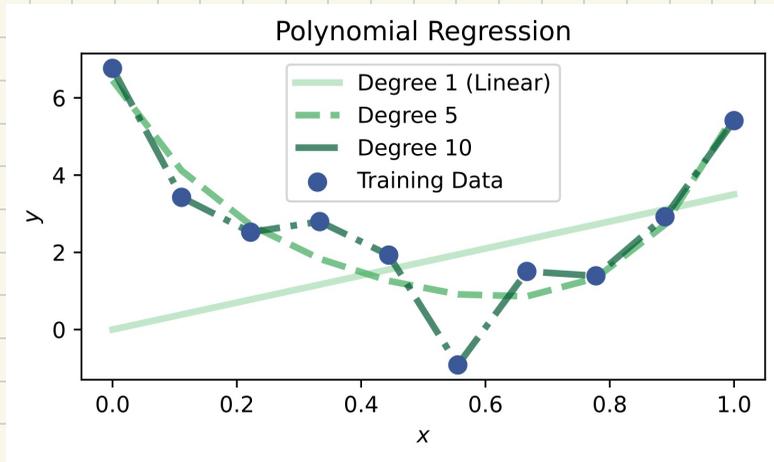
Claim:

$$\min_{w \in \mathbb{R}^{d_{\text{aug}}}} \|y - X_{\text{aug}} w\|_2^2 \leq \min_{w \in \mathbb{R}^d} \|y - X w\|_2^2$$

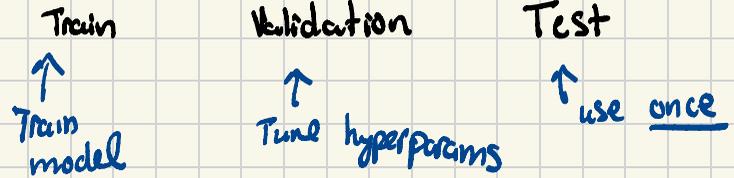
Downsides of more features?

1. Time complexity
(mitigated by kernel methods)

2. Overfitting



Generalization Error



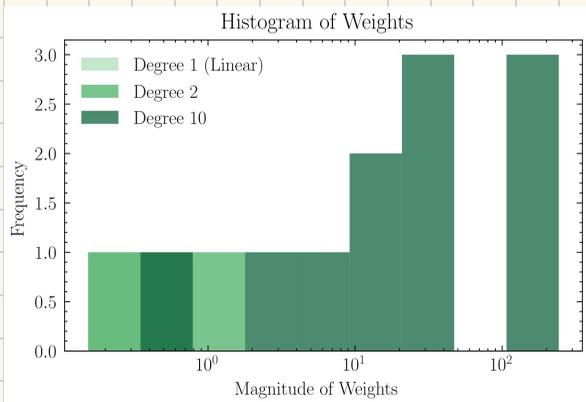
Overfit with :

1. Too many parameters
2. Too much training

Regularization

Occam's razor: Simplest explanation
is usually right

Complicated \approx big changes
 \approx big weights

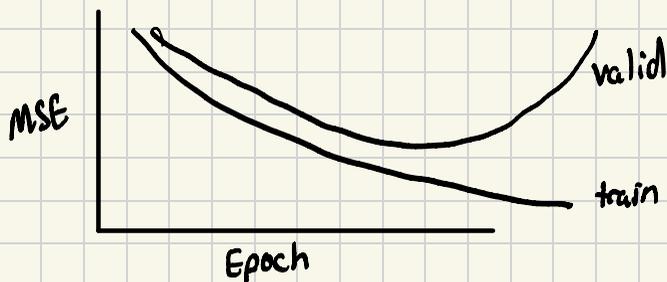


$$\operatorname{argmin}_{\omega}$$

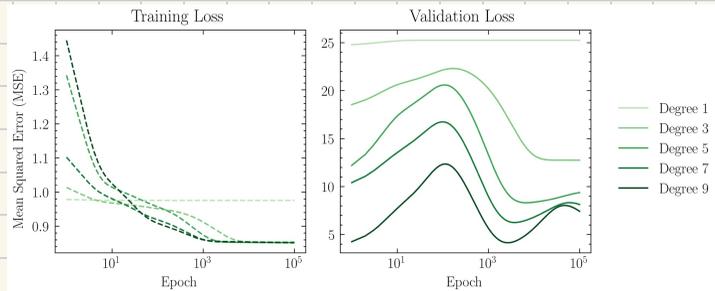
$$\mathcal{L}(\omega) + \lambda \|\omega\|_2^2$$

Traditional View

Tradeoff between generalization
and training



Double Descent

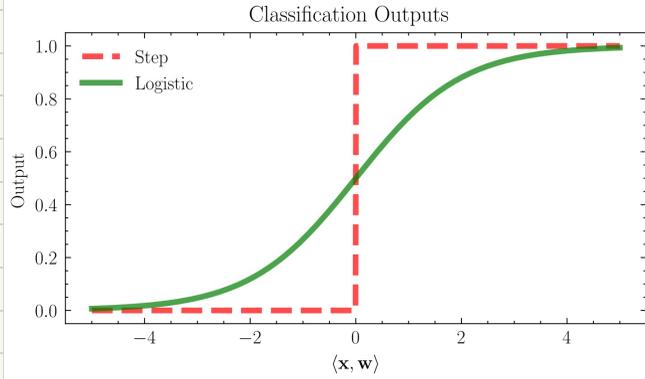


Thursday, February 19

Logistic regression!

Sigmoid $\sigma: \mathbb{R} \rightarrow \mathbb{R}$

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



- smooth
- non-linear
- image is $(0,1)$

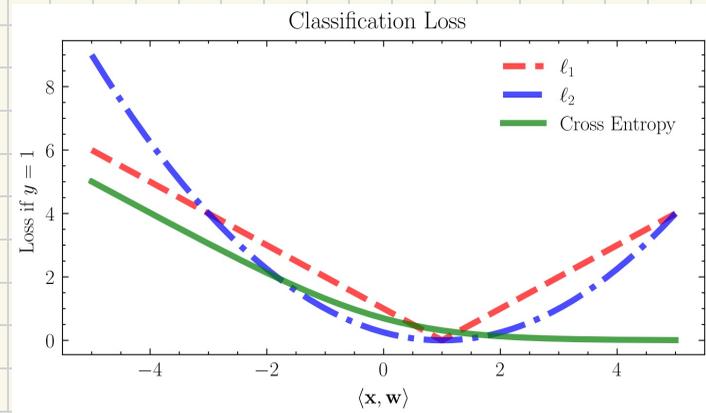
\Rightarrow A probability!! 😊

... and thresholding still works

Loss

$$y \in \{0,1\} \quad f(x) \in (0,1)$$

$$y \log[f(x)] + (1-y) \log[1-f(x)]$$



⊖ If $y = 1$, loss is $\log(f(x))$

😊 If $f(x) \approx 1$, loss is tiny

☹ If $f(x) \approx 0$, loss is huge

Maximum Likelihood Perspective

$$\text{Suppose } y = \begin{cases} 1 & \text{w.p. } f(x) \\ 0 & \text{else} \end{cases}$$

$$\Pr(y \mid x, w) = f(x)^y (1 - f(x))^{1-y}$$

$$\operatorname{argmax}_w \Pr(\text{observing data})$$

$$= \operatorname{argmax}_w \prod_{i=1}^n P(y^{(i)} \mid x^{(i)}, w)$$

= ...

Optimization

exact or gradient descent

Both need the gradient, so...

Lemma: $\nabla_{\omega} \mathcal{L}(\omega) = X^T (\sigma(X\omega) - y)$

$$z = \langle \omega, x \rangle, \quad \hat{y} = \sigma(z)$$

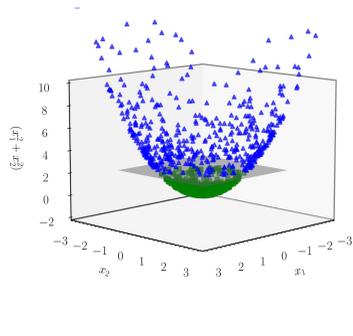
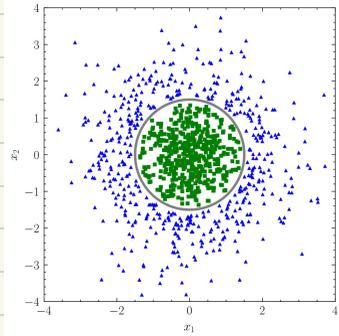
Fact 1: $\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z)) = \hat{y}(1 - \hat{y})$

$$\ell = -y \log \hat{y} - (1-y) \log (1 - \hat{y})$$

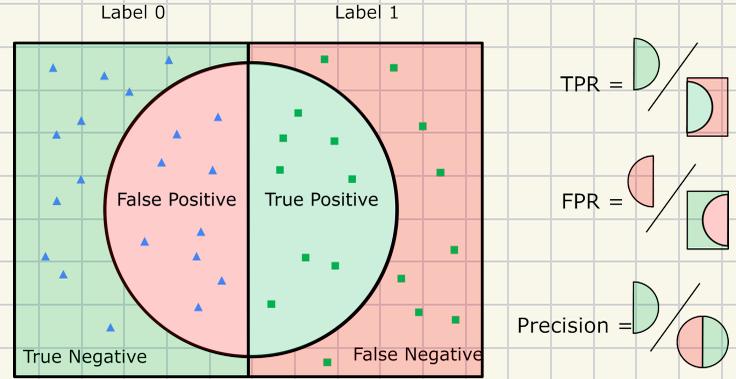
Fact 2: $\frac{\partial \ell}{\partial \hat{y}} = -\frac{y}{\hat{y}} + \frac{(1-y)}{(1-\hat{y})} = \frac{\hat{y} - y}{\hat{y}(1-\hat{y})}$

$$\begin{aligned} \nabla_{\omega} \mathcal{L}(\omega) &= \sum_{i=1}^n \frac{\partial \ell^{(i)}}{\partial \hat{y}^{(i)}} \frac{\partial \hat{y}^{(i)}}{\partial z^{(i)}} \nabla_{\omega} z^{(i)} \\ &= \dots \end{aligned}$$

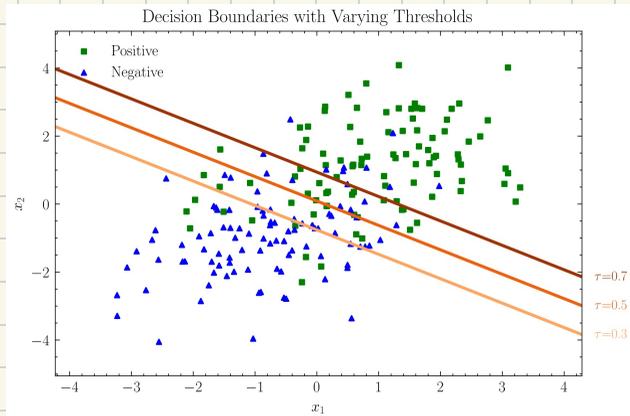
Non-linear Transformations



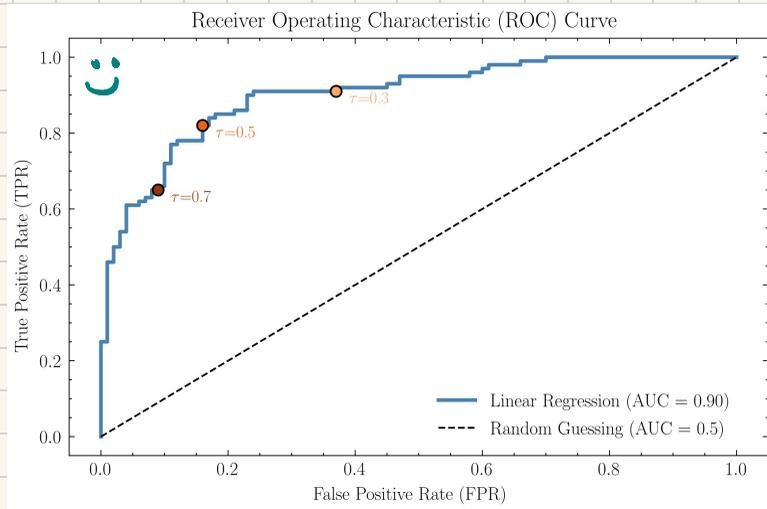
Measuring Error



Varying Threshold



Receiver Operating Characteristic Curve



Output 1 if $\sigma(\langle w, x \rangle) > \frac{1}{2}$

Output 0 else

Decreasing τ :

↑ TPR

↑ FPR