

Machine Learning

4771

Instructor: Tony Jebara

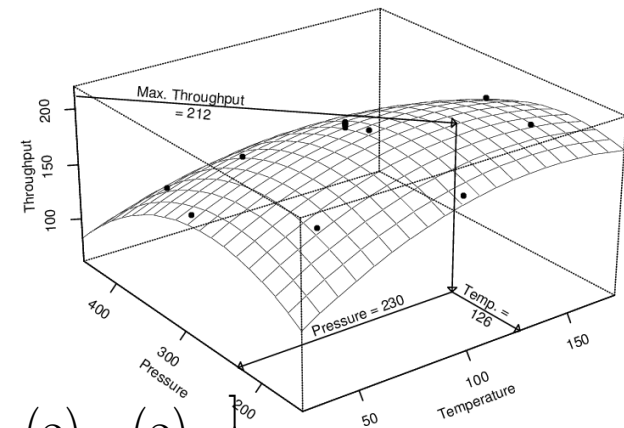
Topic 3

- Additive Models and Linear Regression
- Sinusoids and Radial Basis Functions
- Classification
- Logistic Regression
- Gradient Descent

Polynomial Basis Functions

- To fit a P'th order polynomial function to multivariate data: concatenate columns of all monomials up to power P
- E.g. 2 dimensional data and 2nd order polynomial (quadratic)

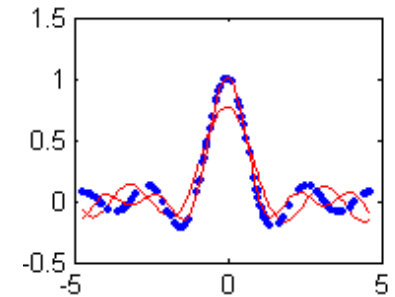
$$\begin{bmatrix} x_1(1) & x_1(2) \\ \vdots & \vdots \\ x_i(1) & x_i(2) \\ \vdots & \vdots \\ x_N(1) & x_N(2) \end{bmatrix} \quad \rightarrow \quad \begin{bmatrix} 1 & x_1(1) & x_1(2) & x_1(1)x_1(1) & x_1(1)x_1(2) & x_1(2)x_1(2) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_i(1) & x_i(2) & x_i(1)x_i(1) & x_i(1)x_i(2) & x_i(2)x_i(2) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_N(1) & x_N(2) & x_N(1)x_N(1) & x_N(1)x_N(2) & x_N(2)x_N(2) \end{bmatrix}$$



Sinusoidal Basis Functions

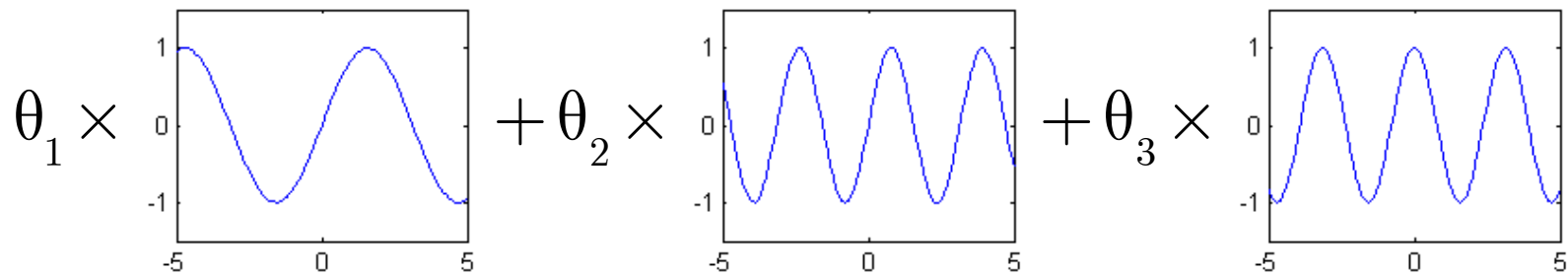
- More generally, we don't just have to deal with polynomials, use any set of basis fn's:

$$f(x; \theta) = \sum_{p=1}^P \theta_p \phi_p(x) + \theta_0$$



- These are generally called **Additive Models**
- Regression adds linear combinations of the basis fn's
- For example: **Fourier (sinusoidal) basis**

$$\phi_{2k}(x_i) = \sin(kx_i) \quad \phi_{2k+1}(x_i) = \cos(kx_i)$$
- Note, don't have to be a basis per se, usually subset



Radial Basis Functions

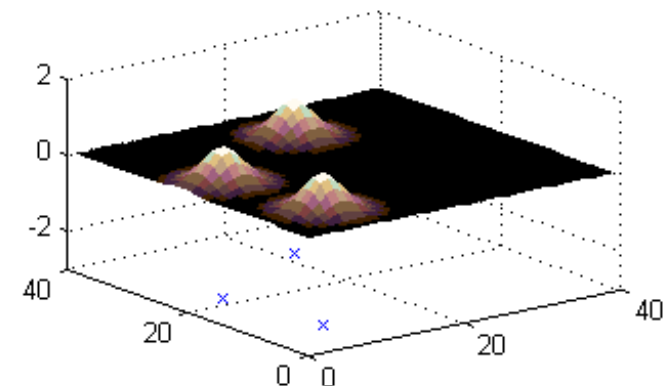
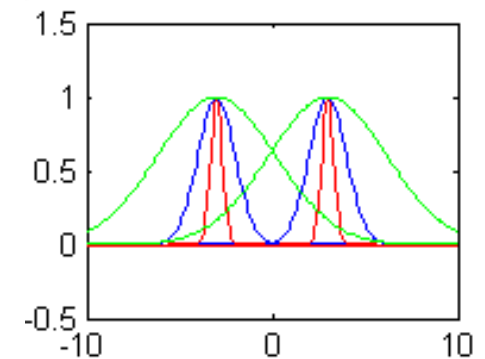
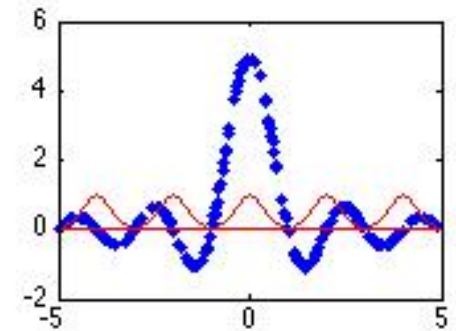
- Can act as prototypes of the data itself

$$f(\mathbf{x}; \theta) = \sum_{k=1}^N \theta_k \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right)$$

- Parameter σ = standard deviation
 σ^2 = covariance

controls how wide bumps are
 what happens if too big/small?

- Also works in multi-dimensions
- Called RBF for short



Radial Basis Functions

- Each training point leads to a bump function

$$f(\mathbf{x}; \theta) = \sum_{k=1}^N \theta_k \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{x}_k\|^2\right)$$

- Reuse solution from linear regression: $\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- Can view the data instead as \mathbf{X} , a big matrix of size $N \times N$

$$\mathbf{X} = \begin{bmatrix} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_1 - \mathbf{x}_1\|^2\right) & \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_1 - \mathbf{x}_2\|^2\right) & \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_1 - \mathbf{x}_3\|^2\right) \\ \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_2 - \mathbf{x}_1\|^2\right) & \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_2 - \mathbf{x}_2\|^2\right) & \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_2 - \mathbf{x}_3\|^2\right) \\ \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_3 - \mathbf{x}_1\|^2\right) & \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_3 - \mathbf{x}_2\|^2\right) & \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_3 - \mathbf{x}_3\|^2\right) \end{bmatrix}$$

- For RBFs, \mathbf{X} is square and symmetric, so solution is just

$$\nabla_{\theta} R = 0 \rightarrow \mathbf{X}^T \mathbf{X} \theta = \mathbf{X}^T \mathbf{y} \rightarrow \mathbf{X} \theta = \mathbf{y} \rightarrow \theta^* = \mathbf{X}^{-1} \mathbf{y}$$

Evaluating Our Learned Function

- We minimized empirical risk to get θ^*
- How well does $f(x; \theta^*)$ perform on future data?
- It should *Generalize* and have low **True Risk**:

$$R_{true}(\theta) = \int P(x, y) L(y, f(x; \theta)) dx dy$$

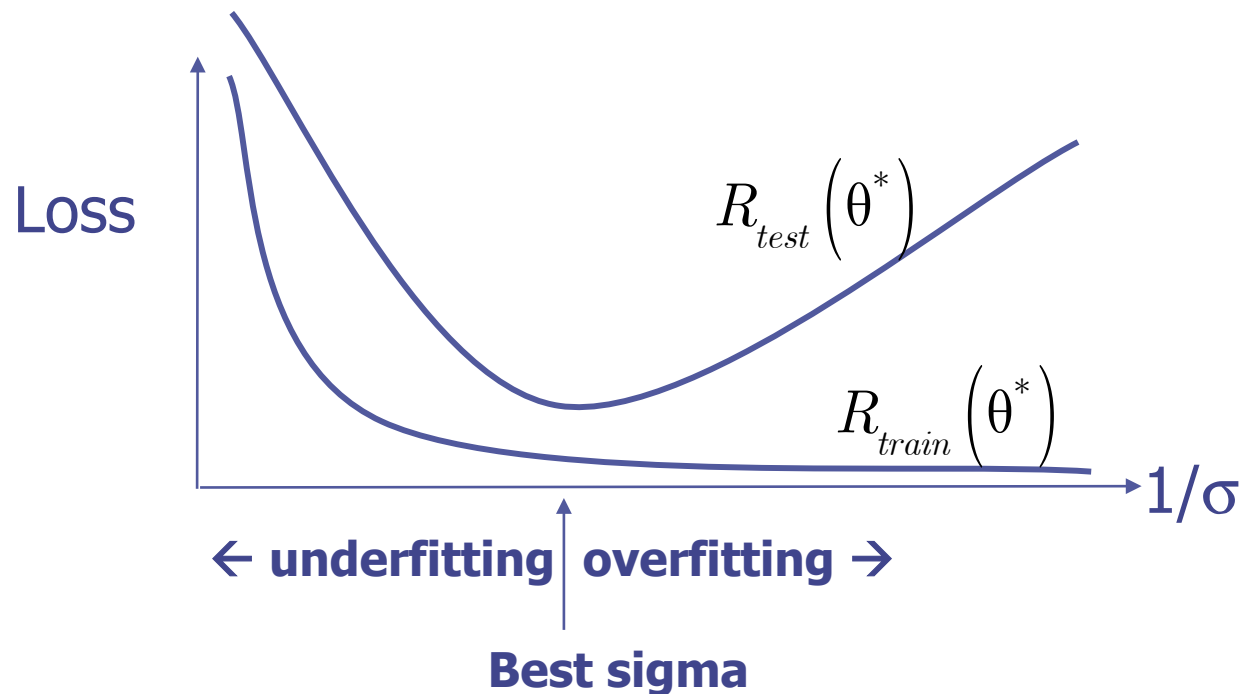
- Can't compute true risk, instead use **Testing Empirical Risk**
- We randomly split data into training and testing portions

$$\left\{ (x_1, y_1), \dots, (x_N, y_N) \right\} \quad \left\{ (x_{N+1}, y_{N+1}), \dots, (x_{N+M}, y_{N+M}) \right\}$$

- Find θ^* with **training data**: $R_{train}(\theta) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i; \theta))$
- Evaluate it with **testing data**: $R_{test}(\theta) = \frac{1}{M} \sum_{i=N+1}^{N+M} L(y_i, f(x_i; \theta))$

Crossvalidation

- Try fitting with different sigma radial basis function widths
- Select sigma which gives lowest $R_{test}(\theta^*)$



- Think of sigma as a measure of the simplicity of the model
- Thinner RBFs are more flexible and complex

Regularized Risk Minimization

- Empirical Risk Minimization gave overfitting & underfitting
- We want to add a penalty for using too many theta values
- This gives us the Regularized Risk

$$\begin{aligned} R_{regularized}(\theta) &= R_{empirical}(\theta) + Penalty(\theta) \\ &= \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i; \theta)) + \frac{\lambda}{2N} \|\theta\|^2 \end{aligned}$$

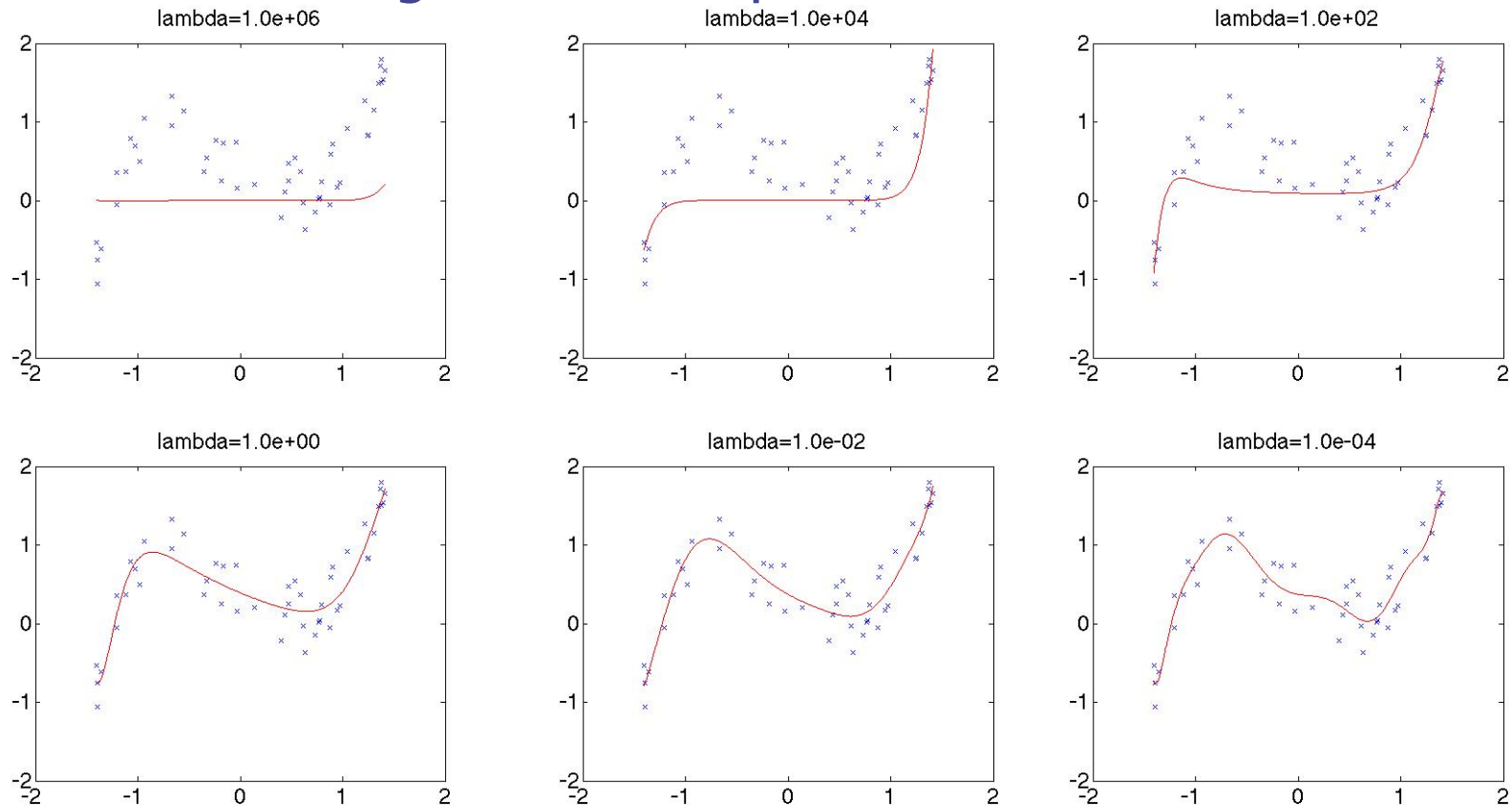
- Solution for Regularized Risk with Least Squares Loss:

$$\nabla_{\theta} R_{regularized} = 0 \quad \Rightarrow \quad \nabla_{\theta} \left(\frac{1}{2N} \|\mathbf{y} - \mathbf{X}\theta\|^2 + \frac{\lambda}{2N} \|\theta\|^2 \right) = 0$$

$$\theta^* = \left(\mathbf{X}^T \mathbf{X} + \lambda I \right)^{-1} \mathbf{X}^T \mathbf{y}$$

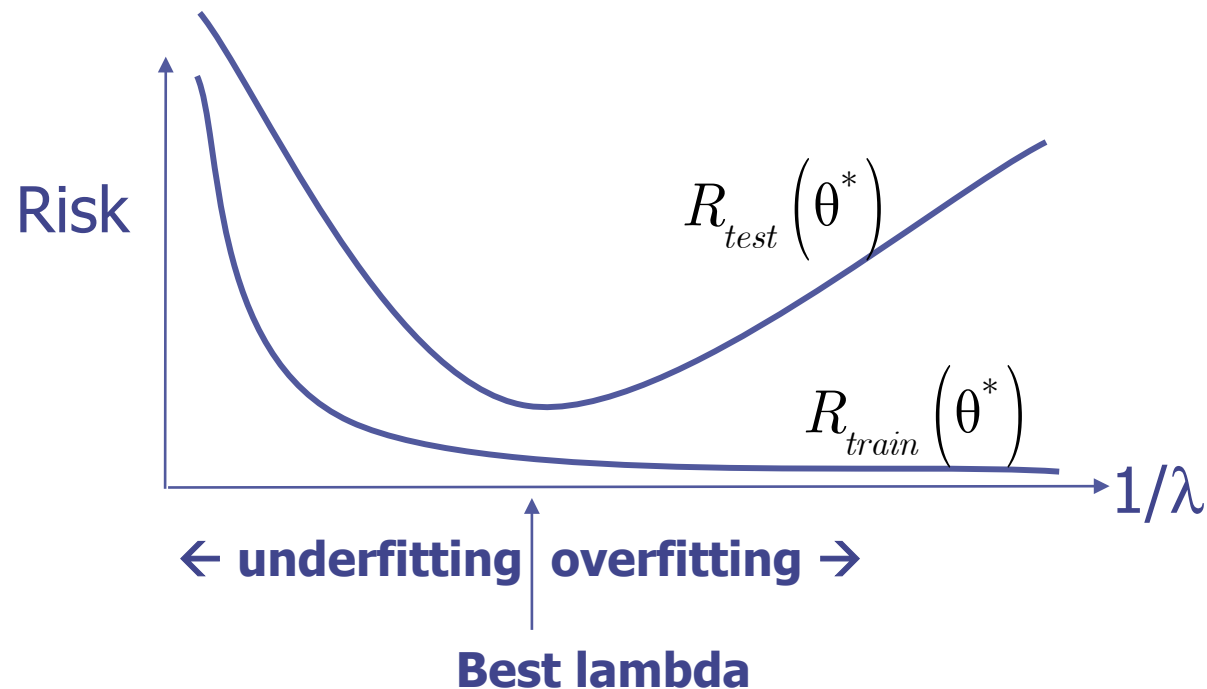
Regularized Risk Minimization

- Have $D=16$ features (or $P=15$ throughout)
- Try minimizing $R_{\text{regularized}}(\theta)$ to get θ^* with different λ
- Note that $\lambda=0$ give back Empirical Risk Minimization



Crossvalidation

- Try fitting with different lambda regularization levels
- Select lambda which gives lowest $R_{test}(\theta^*)$



- Lambda measures simplicity of the model
- Models with low lambda are more flexible

From Regression To Classification

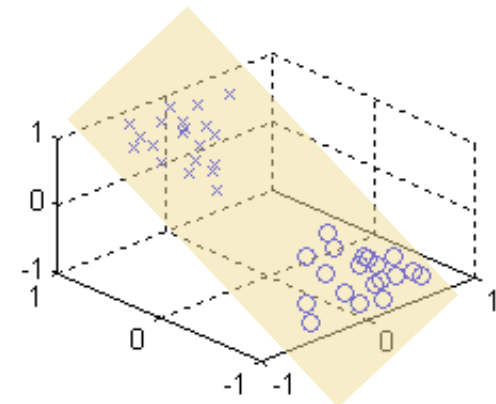
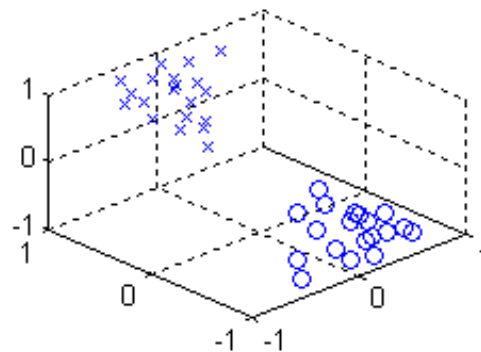
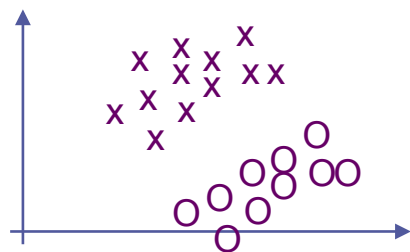
- Classification is another important learning problem

Regression $\mathcal{X} = \left\{ \left(\mathbf{x}_1, y_1 \right), \left(\mathbf{x}_2, y_2 \right), \dots, \left(\mathbf{x}_N, y_N \right) \right\}$ $\mathbf{x} \in \mathbb{R}^D$ $y \in \mathbb{R}^1$

Classification $\mathcal{X} = \left\{ \left(\mathbf{x}_1, y_1 \right), \left(\mathbf{x}_2, y_2 \right), \dots, \left(\mathbf{x}_N, y_N \right) \right\}$ $\mathbf{x} \in \mathbb{R}^D$ $y \in \{0, 1\}$

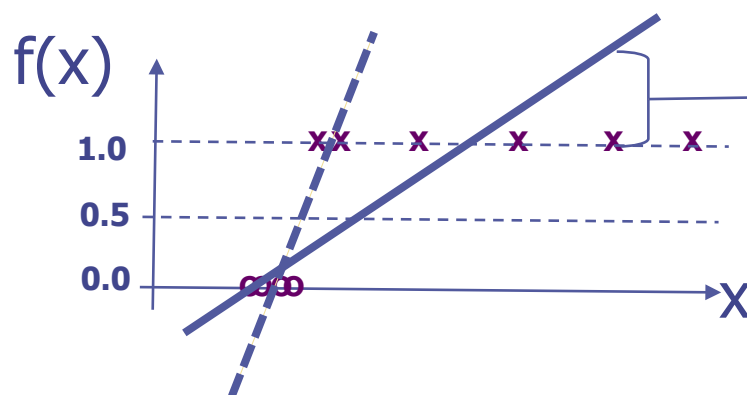
- E.g. Given $\mathbf{x} = [\text{tumor size, tumor density}]$
Predict y in $\{\text{benign, malignant}\}$

- Should we solve this as a least squares regression problem?



Classification vs. Regression

- a) Classification needs binary answers like $\{0,1\}$
 - b) Least squares is an unfair measure of risk here
 - e.g. Why penalize a correct but large positive y answer?
 - e.g. Why penalize a correct but large negative y answer?
- Example: not good to use regression output for a decision
 - $f(x) > 0.5 \rightarrow \text{Class 1}$ $f(x) < 0.5 \rightarrow \text{Class 0}$
 - if $f(x) = -3.8$ & correct class = 0, squared error penalizes it...



We pay a hefty squared error loss here even if we got the correct classification result. The thick solid line model makes two mistakes while the dashed model is perfect

Classification vs. Regression

We will consider the following four steps to improve from naïve regression to get better classification learning:

- 1) Fix functions $f(x)$ to give binary output (logistic neuron)
- 2) Fix our definition of the Risk we will minimize so that we get good classification accuracy (logistic loss)

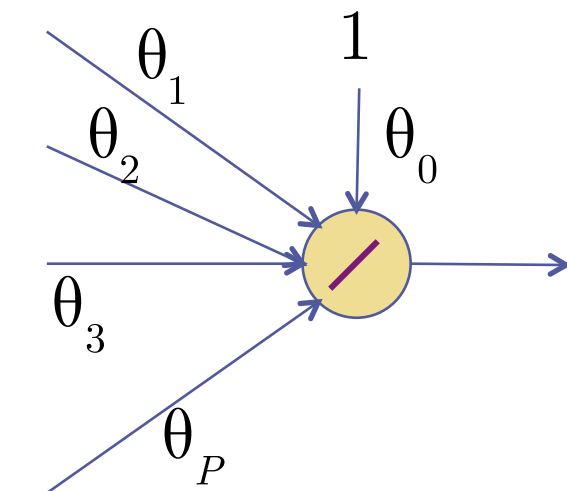
...and later on...

- 3) Make an even better fix on $f(x)$ to binarize (perceptron)
- 4) Make an even better risk (perceptron loss)

Logistic Neuron (McCullough-Pitts)

- To output binary, use squashing function $g()$.

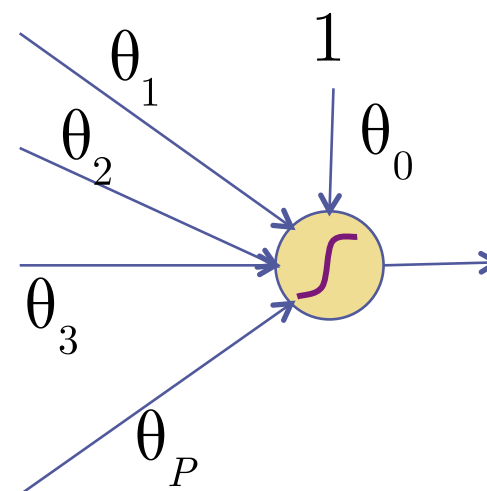
$$f(\mathbf{x}; \theta) = \theta^T \mathbf{x}$$



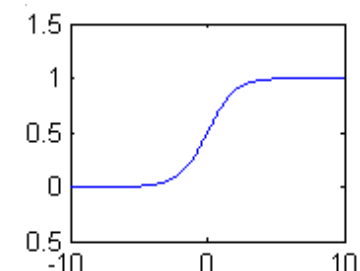
Linear neuron

$$f(\mathbf{x}; \theta) = g(\theta^T \mathbf{x})$$

$$g(z) = \left(1 + \exp(-z)\right)^{-1}$$



Logistic Neuron



- This squashing is called **sigmoid** or **logistic function**

Logistic Regression

- Given a classification problem with binary outputs

$$\mathcal{X} = \left\{ (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N) \right\} \quad \mathbf{x} \in \mathbb{R}^D \quad y \in \{0, 1\}$$

- Use this function and output 1 if $f(\mathbf{x}) > 0.5$ and 0 otherwise

$$f(\mathbf{x}; \theta) = \left(1 + \exp(-\theta^T \mathbf{x}) \right)^{-1}$$

Short hand for Linear Functions

- What happened to adding the intercept?

$$f(\mathbf{x}; \theta) = \theta^T \mathbf{x} + \theta_0$$

$$= \begin{bmatrix} \theta(1) \\ \theta(2) \\ \vdots \\ \theta(D) \end{bmatrix}^T \begin{bmatrix} \mathbf{x}(1) \\ \mathbf{x}(2) \\ \vdots \\ \mathbf{x}(D) \end{bmatrix} + \theta_0 = \begin{bmatrix} \theta_0 \\ \theta(1) \\ \theta(2) \\ \vdots \\ \theta(D) \end{bmatrix}^T \begin{bmatrix} 1 \\ \mathbf{x}(1) \\ \mathbf{x}(2) \\ \vdots \\ \mathbf{x}(D) \end{bmatrix} = \vec{\theta}^T \vec{\mathbf{x}}$$

Logistic Regression

- Given a classification problem with binary outputs

$$\mathcal{X} = \left\{ (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N) \right\} \quad \mathbf{x} \in \mathbb{R}^D \quad y \in \{0, 1\}$$

- Fix#1: use $f(\mathbf{x})$ below, output 1 if $f(\mathbf{x}) > 0.5$ and 0 otherwise

$$f(\mathbf{x}; \theta) = \left(1 + \exp(-\theta^T \mathbf{x}) \right)^{-1}$$

Logistic Regression

- Given a classification problem with binary outputs

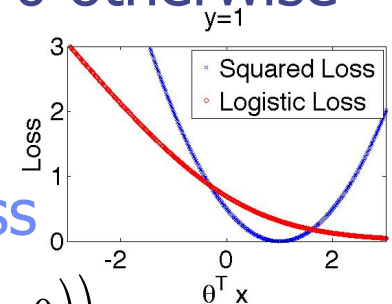
$$\mathcal{X} = \left\{ (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N) \right\} \quad \mathbf{x} \in \mathbb{R}^D \quad y \in \{0, 1\}$$

- Fix#1: use $f(\mathbf{x})$ below, output 1 if $f(\mathbf{x}) > 0.5$ and 0 otherwise

$$f(\mathbf{x}; \theta) = \left(1 + \exp(-\theta^T \mathbf{x}) \right)^{-1}$$

- Fix#2: instead of squared loss, use **Logistic Loss**

$$L_{\log} \left(y, f(\mathbf{x}; \theta) \right) = (y_i - 1) \log(1 - f(\mathbf{x}; \theta)) - y_i \log(f(\mathbf{x}; \theta))$$



- This method is called **Logistic Regression**.
- But Empirical Risk Minimization has no closed-form sol'n:

$$R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N \left(y_i - 1 \right) \log \left(1 - f(\mathbf{x}_i; \theta) \right) - y_i \log \left(f(\mathbf{x}_i; \theta) \right)$$

Logistic Regression

- With logistic squashing function, minimizing $R(\theta)$ is harder

$$R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - 1) \log(1 - f(\mathbf{x}_i; \theta)) - y_i \log(f(\mathbf{x}_i; \theta))$$

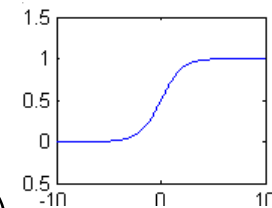
$$\nabla_{\theta} R = \frac{1}{N} \sum_{i=1}^N \left(\frac{1 - y_i}{1 - f(\mathbf{x}_i; \theta)} - \frac{y_i}{f(\mathbf{x}_i; \theta)} \right) f'(\mathbf{x}_i; \theta) = 0 \quad ???$$

- Can't minimize risk and find best theta analytically!
- Let's try finding best theta numerically.
- Use the following to compute gradient

$$f(\mathbf{x}; \theta) = \left(1 + \exp(-\theta^T \mathbf{x}) \right)^{-1} = g(\theta^T \mathbf{x})$$

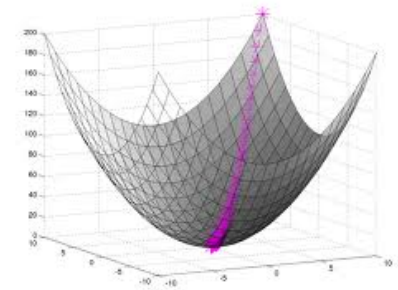
- Here, $g()$ is the logistic squashing function

$$g(z) = \left(1 + \exp(-z) \right)^{-1} \quad g'(z) = g(z)(1 - g(z))$$



Gradient Descent

- Useful when we can't get minimum solution in closed form
- Gradient points in direction of fastest increase
- Take step in the opposite direction!
- Gradient Descent Algorithm



choose scalar step size η , & tolerance ϵ
initialize $\theta^0 = \text{small random vector}$

$$\theta^1 = \theta^0 - \eta \nabla_{\theta} R_{emp} \Big|_{\theta^0}, \quad t = 1$$

$$\text{while } \left\| \theta^t - \theta^{t-1} \right\| \geq \epsilon \quad \left\{ \begin{array}{l} \theta^{t+1} = \theta^t - \eta \nabla_{\theta} R_{emp} \Big|_{\theta^t}, \quad t = t + 1 \end{array} \right. \quad \left. \right\}$$

- For appropriate η , this will converge to local minimum

Logistic Regression

- Logistic regression gives better classification performance
- Its empirical risk is

$$R_{emp}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - 1) \log(1 - f(\mathbf{x}_i; \theta)) - y_i \log(f(\mathbf{x}_i; \theta))$$

- This $R(\theta)$ is convex so gradient descent always converges to the same solution

- Make predictions using

$$f(\mathbf{x}; \theta) = \left(1 + \exp(-\theta^T \mathbf{x})\right)^{-1}$$

- Output 1 if $f > 0.5$
- Output 0 otherwise

