

Supervised machine learning

- Machine learning: methods that improve with access to more data.

Examples in genomics:

- Identify sequence elements: Genes, exons, splice sites, enhancers, positioned nucleosomes.
 - Assign functional annotations to genes: GO terms, co-functionality.
 - Dynamics: Gene expression, regulatory relationships, structure
 - Input: sequence, functional genomics, ~~the~~ gene expression, ~~sequences~~ populations of sequences
- Blurry line between ML and statistics. ML focuses on large, heterogeneous data sets.
 - Supervised ML: predict "label". (distinct from unsupervised)
 - Appropriate when prediction > interpretation
 - ~~the~~ Prediction is the real task: Diagnosis, drug side effects, DNA primer will bind
 - Is it predictable? Gene expression, splice sites

Problem setup

- Example: Predict gene expression from functional genomics data.

- ML setup:

$$x: \boxed{x_1 \mid x_2 \mid \dots}$$

"Label" — prediction target (y)

"features" — anything used for prediction ($x_{1..n}$)

- Assume label-feature pairs are independent draws from the same distribution.

- Gene expression example:

y : is gene expressed? (RNA-seq > cutoff)

x_1 : is H3K36me3 present at promoter? (reads > cutoff)

- Classification: 0/1 labels

Regression: real-valued labels

Logistic regression

θ_j : Association of x_j with y

$$x^{(i)}: \begin{array}{|c|c|c|} \hline x_1 & x_2 & \dots \\ \hline \end{array}$$

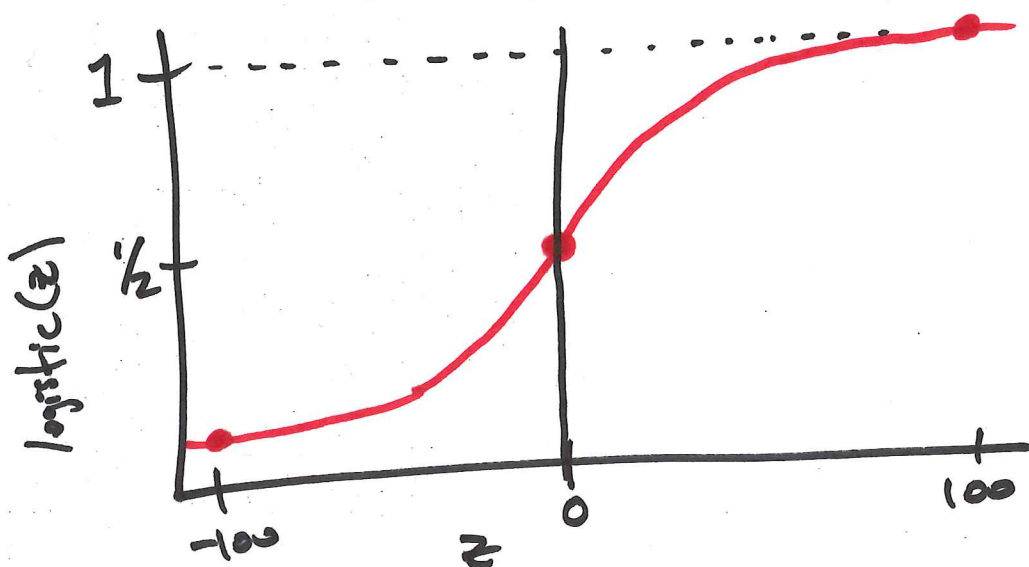
confidence that $y^{(i)} = 1$:

$$\theta: \begin{array}{|c|c|c|} \hline \theta_1 & \theta_2 & \dots \\ \hline \end{array}$$

$$\sum_i x_i \cdot \theta_i$$

Map confidence \rightarrow probability: logistic function

$$\text{logistic}(z) = \frac{\exp(z)}{1 + \exp(z)}$$

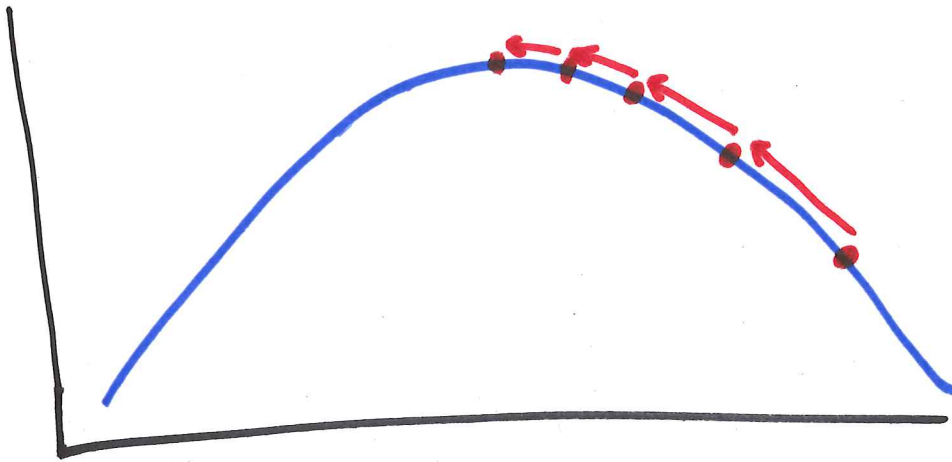


$$h_{\theta}(x) = \text{logistic}\left(\sum_j x_j \theta_j\right)$$

$$P_{\theta}(y|x) = \prod_i P_{\theta}(y_i|x_i)$$

Optimization

- General problem: Find the maximum of an "objective function".
- Gradient descent



- Derivative of logistic function

$$\text{logistic}(z) = \frac{\exp(z)}{1 + \exp(z)} = \frac{1}{\exp(-z) + 1}$$

$$\begin{aligned} \frac{d(\text{logistic}(z))}{dz} &= \left(\frac{1}{\exp(-z) + 1} \right)^2 \exp(-z) \\ &= \text{logistic}(z) \frac{\exp(-z) + 1 - 1}{\exp(-z) + 1} \end{aligned}$$

$$= \text{logistic}(z) (1 - \text{logistic}(z))$$

Optimization 2

$$h_{\theta}(x) = \text{logistic}\left(\sum_j \theta_j x_j\right)$$

~~$$P(y^{(i)} | x^{(i)}) = h_{\theta}(x^{(i)})^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{(1-y^{(i)})}$$~~

$$\ell(\theta) = \log(P_{\theta}(y|x)) = \sum_i y \log h_{\theta}(x) + (1-y) \log(1 - h_{\theta}(x))$$

$$\frac{d\ell(\theta)}{d\theta_j} = \sum_i \left[y^{(i)} \frac{1}{h_{\theta}(x^{(i)})} - (1-y^{(i)}) \frac{1}{1-h_{\theta}(x^{(i)})} \right] \left(\frac{d}{d\theta_j} h_{\theta}(x^{(i)}) \right)$$

$$= \sum_i \left[\right] h_{\theta}(x)(1-h_{\theta}(x)) \left(\frac{d}{d\theta_j} \sum_{j'} x_{j'} \theta_{j'} \right)$$

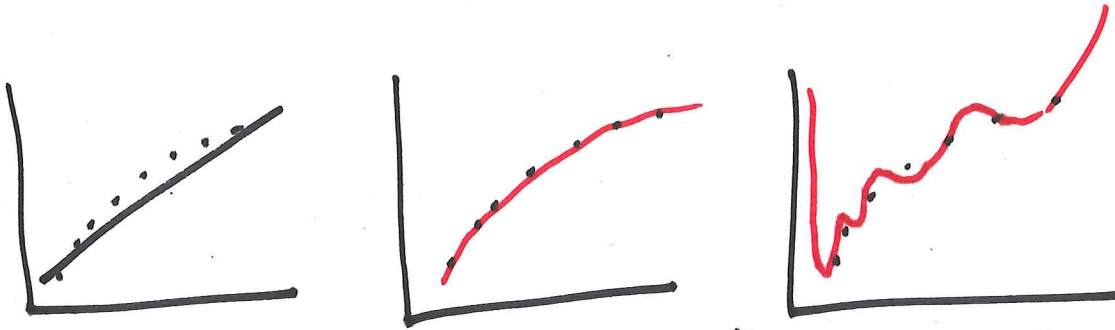
$$= \sum_i \left[y(1-h_{\theta}(x)) - (1-y)h_{\theta}(x) \right] x_j$$

$$= \sum_i \left[y - h_{\theta}(x) \right] x_j$$

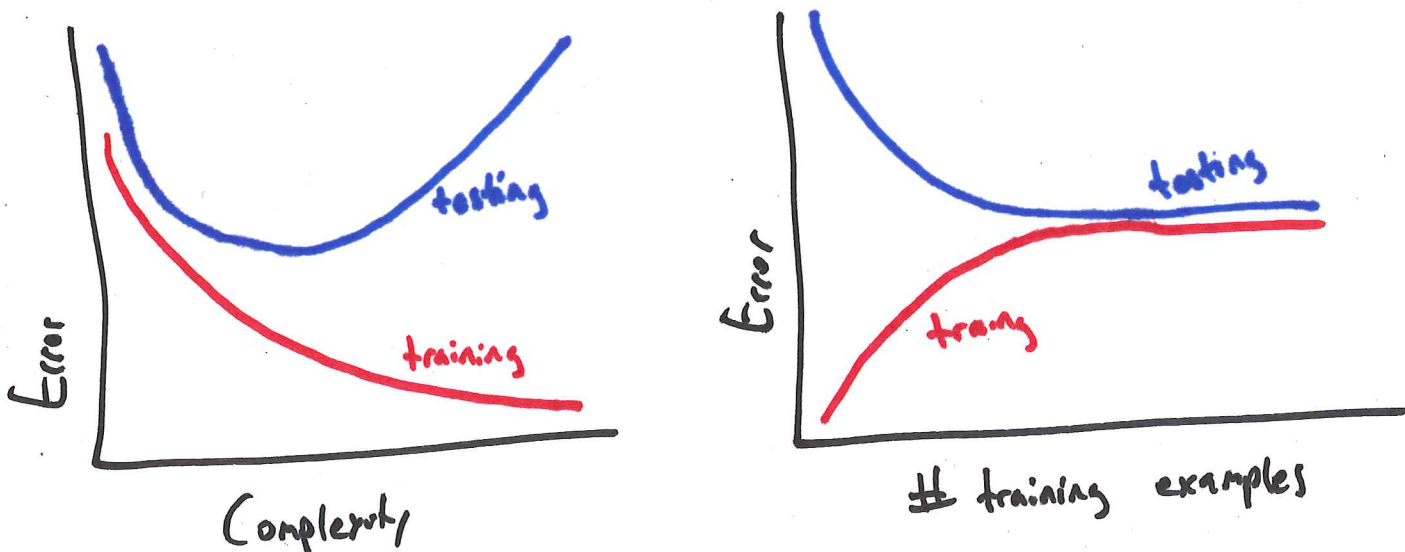
$$\theta_j := \theta_j + \alpha \sum_i (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Overfitting and model complexity

- Example: fitting a polynomial



- How well does a model ~~generalize~~ ^{generalize} to data not in training set?
- Tradeoff: model complexity



- Use this behavior to guide your choice of models:
 - features ; ML algorithm ; regularization/prior