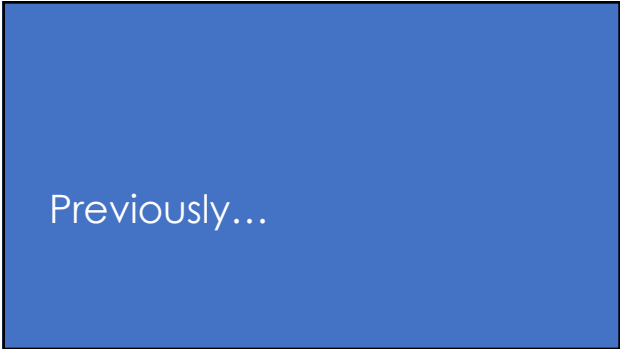
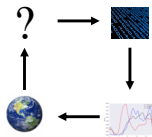


Classification & Logistic Regression & maybe deep learning

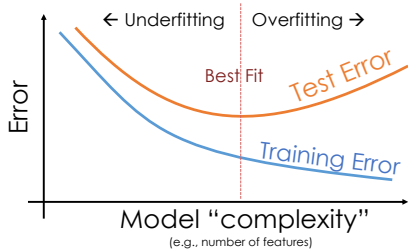
Slides by:

Joseph E. Gonzalez jgonzal@cs.berkeley.edu



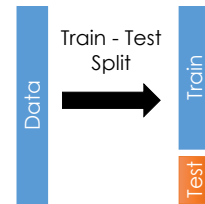
Training vs Test Error

Training error typically underestimates test error.



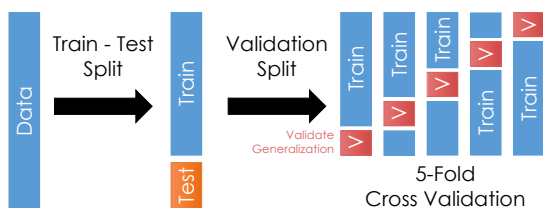
Generalization: The Train-Test Split

- > **Training Data:** used to fit model
- > **Test Data:** check generalization error
- > How to split?
 - > Randomly, Temporally, Geo...
 - > Depends on application (usually randomly)
- > What size? (90%-10%)
 - > Larger training set → more complex models
 - > Larger test set → better estimate of generalization error
 - > Typically between 75%-25% and 90%-10%



You can only use the test dataset once after deciding on the model.

Generalization: Validation Split



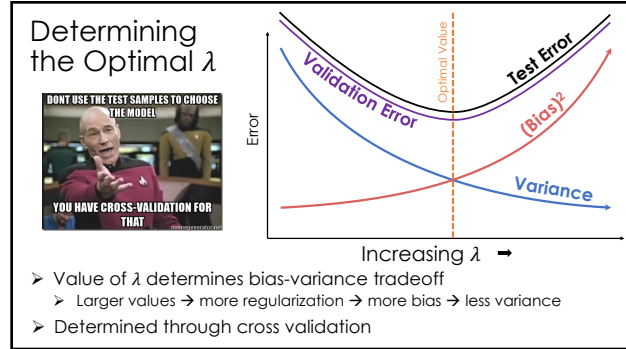
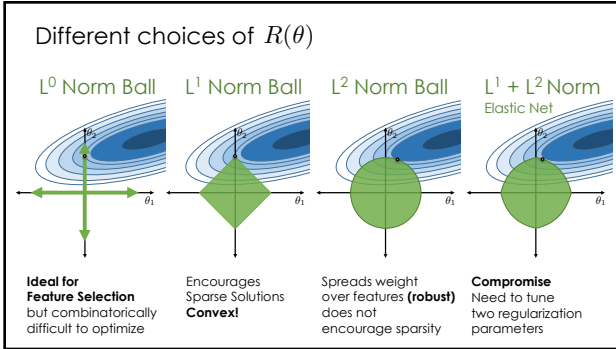
Cross validation simulates multiple train test-splits on the training data.

Regularized Loss Minimization

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, f_{\theta}(x_i)) + \lambda R(\theta)$$

Regularization Parameter

- > Larger values of λ → more regularization
- > **Confusing!**: Scikit-learn uses $\alpha = 1/\lambda$
 - > Larger values of α → less regularization.



Using Scikit-Learn for Regularized Regression

```
import sklearn.linear_model
```

- Confusion Warning:** Regularization parameter $\alpha = 1/\lambda$
 - larger $\alpha \rightarrow$ less regularization \rightarrow greater complexity \rightarrow overfitting
- Lasso Regression (L1)**
 - `linear_model.Lasso(alpha=3.0)`
 - `linear_model.LassoCV()` automatically picks α by cross-validation
- Ridge Regression (L2)**
 - `linear_model.Ridge(alpha=3.0)`
 - `linear_model.RidgeCV()` automatically selects α by cross-validation
- Elastic Net (L1 + L2)**
 - `linear_model.ElasticNet(alpha=3.0, l1_ratio = 2.0)`
 - `linear_model.ElasticNetCV()` automatically picks α by cross-validation

Standardization and the Intercept Term

Height = $\theta_1 \text{age_in_seconds} + \theta_2 \text{weight_in_tons}$

Small, Large

- Regularization penalized dimensions equally
- Standardization**
 - Ensure that each dimensions has the same scale
 - centered around zero
- Intercept Terms**
 - Typically don't regularize intercept term
 - Center y values (e.g., subtract mean)

Standardization

For each dimension k:

$$z_k = \frac{x_k - \mu_k}{\sigma_k}$$

Regularization and High-Dimensional Data

Regularization is often used with high-dimensional data

d

n

Φ

Tall Skinny Matrix

- $n \gg d$
- typically dense
- Regularization can help with complex feature transformations

d

n

Φ

High-dimensional sparse matrix

- $d > n$
- requires regularization
- Goal: to determine **informative dimensions**
 - Consider L1 (Lasso) Regularization.
- Goal: is to make **robust predictions**
 - Consider L2 (+L1) Regularization

Today Classification

So far

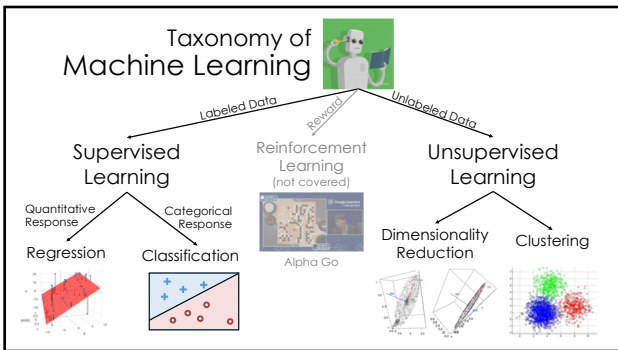
$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2 + \lambda \mathbf{R}(\theta)$$

Squared Loss

Regularization

Classification

isAlive?
{0,1}



- ### Kinds of Classification
- Predicting a categorical variable
- > **Binary** classification: Two classes
 - > Examples: Spam/Not Spam, churn/stay
 - > **Multiclass** classification: Many classes (>2)
 - > Examples: Image labeling (Cat, Dog, Car), Next word in a sentence ...
 - > **Structured prediction** tasks (Classification)
 - > Multiple related predictions
 - > Examples: Translation, Voice recognition

Classification

isAlive?
{0,1}

Can we just use least squares?

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2 + \lambda \mathbf{R}(\theta)$$

Squared Loss

Python Demo

Classification

Can we just use least squares?

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2 + \lambda R(\theta)$$

Classification

Can we just use least squares?

- > Yes ...
- > Need ...
- > **Don't use Least Squares for Classification**
- > Interpret model ...
- > Sensitive to outliers

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2 + \lambda R(\theta)$$

Defining a New Model for Classification

Logistic Regression

- > Widely used models for **binary classification**:

$x = \text{"Get a FREE sample ..."} \Rightarrow y = 1$

$\phi(x) = [2.0, 0.0, \dots, 1.0, 0.5]$

$1 = \text{"Spam"}$
 $0 = \text{"Ham"}$

Why is ham good and spam bad? ...
<https://www.youtube.com/watch?v=cmwv2M138E>

- > Models the probability of $y=1$ given x

$$\hat{P}_{\theta}(y = 1 | x) = \sigma(\phi(x)^T \theta) = \frac{1}{1 + \exp(-\phi(x)^T \theta)}$$

Logistic Regression

Model: $\hat{P}_{\theta}(y = 1 | x) = \sigma(\phi(x)^T \theta) = \frac{1}{1 + \exp(-\phi(x)^T \theta)}$

Generalized Linear Model:
Non-linear transformation of a linear model.

- > Widely used models for **binary classification**:

$x = \text{"Get a FREE sample ..."} \Rightarrow y = 1$

$\phi(x) = [2.0, 0.0, \dots, 1.0, 0.5]$

$1 = \text{"Spam"}$
 $0 = \text{"Ham"}$

Why is ham good and spam bad? ...
<https://www.youtube.com/watch?v=cmwv2M138E>

- > Models the probability of $y=1$ given x

$$\hat{P}_{\theta}(y = 1 | x) = \sigma(\phi(x)^T \theta) = \frac{1}{1 + \exp(-\phi(x)^T \theta)}$$

$$\hat{P}_{\theta}(y = 0 | x) = 1 - \hat{P}_{\theta}(y = 1 | x)$$

Python Demo

The Logistic Regression Model

$$\text{Model: } \hat{\mathbf{P}}_{\theta}(y = 1 | x) = \sigma(\phi(x)^T \theta) = \frac{1}{1 + \exp(-\phi(x)^T \theta)}$$

How do we fit the model to the data?

Defining the Loss

Could we use the Squared Loss

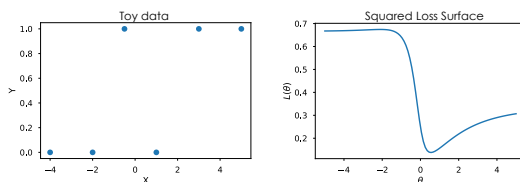
➤ What about squared loss and the new model:

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\phi(x_i)^T \theta))^2$$

- Tries to match probability with 0/1 labels.
- Occasionally used in some neural network applications
- **Non-convex!**

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\phi(x_i)^T \theta))^2$$

- Tries to match probability with 0/1 labels.
- Occasionally used in some neural network applications
- **Non-convex!**



Defining the Cross Entropy Loss

Loss Function

- We want our model to be close to the data:

$$\hat{\mathbf{P}}_{\theta}(y = 1 | x) \approx \mathbf{P}(y = 1 | x)$$

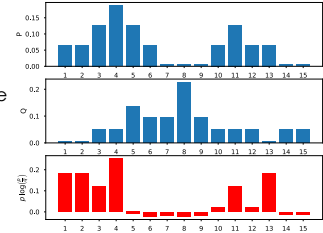
- Kullback–Leibler (KL) Divergence provides a measure of similarity between two distributions:
 - Between two discrete distributions P and Q

$$\mathbf{D}(P||Q) = \sum_{k=1}^K P(k) \log \left(\frac{P(k)}{Q(k)} \right)$$

Kullback–Leibler (KL) Divergence

$$\mathbf{D}(P||Q) = \sum_{k=1}^K P(k) \log \left(\frac{P(k)}{Q(k)} \right)$$

- The average log difference between P and Q weighted by P
- Does not penalize mismatch for rare events with respect to P



Loss Function

- We want our model to be close to the data:

$$\hat{\mathbf{P}}_{\theta}(y = 1 | x) \approx \mathbf{P}(y = 1 | x)$$

- Kullback–Leibler (KL) divergence for classification
 - For a **single** (x,y) data point

$$\mathbf{D}_{KL}(\mathbf{P} || \hat{\mathbf{P}}_{\theta}) = \sum_{k=1}^{K \neq 2} \mathbf{P}(y = k | x) \log \left(\frac{\mathbf{P}(y = k | x)}{\hat{\mathbf{P}}_{\theta}(y = k | x)} \right)$$

- Average KL Divergence for all the data:

- Kullback–Leibler (KL) divergence for classification
 - For a **single** (x,y) data point

$$\mathbf{D}_{KL}(\mathbf{P} || \hat{\mathbf{P}}_{\theta}) = \sum_{k=1}^{K \neq 2} \mathbf{P}(y = k | x) \log \left(\frac{\mathbf{P}(y = k | x)}{\hat{\mathbf{P}}_{\theta}(y = k | x)} \right)$$

- Average KL Divergence for all the data:

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K \mathbf{P}(y_i = k | x_i) \log \left(\frac{\mathbf{P}(y_i = k | x_i)}{\hat{\mathbf{P}}_{\theta}(y_i = k | x_i)} \right)$$

Doesn't depend on θ

$$\mathbf{P}(y_i = k | x_i) \log(\mathbf{P}(y_i = k | x_i)) - \mathbf{P}(y_i = k | x_i) \log(\hat{\mathbf{P}}_{\theta}(y_i = k | x_i))$$

- Average **cross entropy loss**

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K - \mathbf{P}(y_i = k | x_i) \log(\hat{\mathbf{P}}_{\theta}(y_i = k | x_i))$$

Summing from k = 0 to 1 and not k = 1 to 2 (to be consistent with 0/1 labels):

$$\arg \min_{\theta} - \frac{1}{n} \sum_{i=1}^n \left[\mathbf{P}(y_i = 0 | x_i) \log(\hat{\mathbf{P}}_{\theta}(y_i = 0 | x_i)) + \mathbf{P}(y_i = 1 | x_i) \log(\hat{\mathbf{P}}_{\theta}(y_i = 1 | x_i)) \right]$$

$\mathbf{P}(y_i = 1 x_i) = y_i$	$\hat{\mathbf{P}}_{\theta}(y_i = 1 x_i) = \sigma(\phi(x_i)^T \theta)$
$\mathbf{P}(y_i = 0 x_i) = (1 - y_i)$	$\hat{\mathbf{P}}_{\theta}(y_i = 0 x_i) = 1 - \sigma(\phi(x_i)^T \theta)$

- Average **cross entropy loss**

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^{K \neq 2} - \mathbf{P}(y_i = k | x_i) \log(\hat{\mathbf{P}}_{\theta}(y_i = k | x_i))$$

$$\arg \min_{\theta} - \frac{1}{n} \sum_{i=1}^n \left[(1 - y_i) \log(1 - \sigma(\phi(x_i)^T \theta)) + y_i \log(\sigma(\phi(x_i)^T \theta)) \right]$$

Rewriting on one line:

$$\arg \min_{\theta} - \frac{1}{n} \sum_{i=1}^n (y_i \log(\sigma(\phi(x_i)^T \theta)) + (1 - y_i) \log(1 - \sigma(\phi(x_i)^T \theta)))$$

➤ Average **cross entropy loss**

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K -\mathbf{P}(y_i = k | x_i) \log(\hat{\mathbf{P}}_{\theta}(y_i = k | x_i))$$

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \log(\sigma(\phi(x_i)^T \theta)) + (1 - y_i) \log(1 - \sigma(\phi(x_i)^T \theta)))$$

Expanding

$$\log(1 - \sigma(\phi(x_i)^T \theta)) - y_i \log(1 - \sigma(\phi(x_i)^T \theta))$$

Grouping Terms Copy & Paste

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n \left(y_i \log\left(\frac{\sigma(\phi(x_i)^T \theta)}{1 - \sigma(\phi(x_i)^T \theta)}\right) + \log(1 - \sigma(\phi(x_i)^T \theta)) \right)$$

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n \left(y_i \log\left(\frac{\sigma(\phi(x_i)^T \theta)}{1 - \sigma(\phi(x_i)^T \theta)}\right) + \log(1 - \sigma(\phi(x_i)^T \theta)) \right)$$

$\sigma(\phi(x)^T \theta) = \frac{1}{1 + \exp(-\phi(x)^T \theta)}$

$$\log \frac{1}{1 - \frac{1}{1 + \exp(-\phi(x_i)^T \theta)} \times (1 + \exp(-\phi(x_i)^T \theta))}$$

$$\stackrel{\text{Alg.}}{=} \log \frac{1}{1 + \exp(-\phi(x_i)^T \theta) - 1} \stackrel{\text{Alg.}}{=} \log \exp(\phi(x_i)^T \theta) \stackrel{\text{Alg.}}{=} \phi(x_i)^T \theta$$

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n \left(y_i \log\left(\frac{\sigma(\phi(x_i)^T \theta)}{1 - \sigma(\phi(x_i)^T \theta)}\right) + \log(1 - \sigma(\phi(x_i)^T \theta)) \right)$$

$\sigma(\phi(x)^T \theta) = \frac{1}{1 + \exp(-\phi(x)^T \theta)}$

= $\phi(x_i)^T \theta$

A Linear mode of the "Log odds"

$$\log\left(\frac{\sigma(\phi(x_i)^T \theta)}{1 - \sigma(\phi(x_i)^T \theta)}\right) = \phi(x_i)^T \theta = \log\left(\frac{\hat{\mathbf{P}}_{\theta}(y_i = 1 | x_i)}{\hat{\mathbf{P}}_{\theta}(y_i = 0 | x_i)}\right)$$

Linear Model Log odds

A Linear mode of the "Log odds"

$$\log\left(\frac{\sigma(\phi(x_i)^T \theta)}{1 - \sigma(\phi(x_i)^T \theta)}\right) = \phi(x_i)^T \theta = \log\left(\frac{\hat{\mathbf{P}}_{\theta}(y_i = 1 | x_i)}{\hat{\mathbf{P}}_{\theta}(y_i = 0 | x_i)}\right)$$

Linear Model Log odds

Implications?

$$\phi(x_i)^T \theta = 0 \Rightarrow \hat{\mathbf{P}}_{\theta}(y_i = 1 | x_i) = \hat{\mathbf{P}}_{\theta}(y_i = 0 | x_i)$$

$$\phi(x_i)^T \theta > 0 \Rightarrow \hat{\mathbf{P}}_{\theta}(y_i = 1 | x_i) > \hat{\mathbf{P}}_{\theta}(y_i = 0 | x_i)$$

$$\phi(x_i)^T \theta < 0 \Rightarrow \hat{\mathbf{P}}_{\theta}(y_i = 1 | x_i) < \hat{\mathbf{P}}_{\theta}(y_i = 0 | x_i)$$

for any positive ϵ

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n \left(y_i \log\left(\frac{\sigma(\phi(x_i)^T \theta)}{1 - \sigma(\phi(x_i)^T \theta)}\right) + \log(1 - \sigma(\phi(x_i)^T \theta)) \right)$$

$\sigma(\phi(x)^T \theta) = \frac{1}{1 + \exp(-\phi(x)^T \theta)}$

= $\phi(x_i)^T \theta$

Substituting the above result:

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(1 - \sigma(\phi(x_i)^T \theta)))$$

$$1 - \frac{1}{1 + \exp(-\phi(x_i)^T \theta)} \stackrel{\text{Defn. of } \sigma}{=} \frac{\exp(\phi(x_i)^T \theta)}{1 + \exp(\phi(x_i)^T \theta)} \stackrel{\text{Alg.}}{=} \frac{1}{1 + \exp(\phi(x_i)^T \theta)}$$

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(1 - \sigma(\phi(x_i)^T \theta)))$$

$$1 - \frac{1}{1 + \exp(-\phi(x_i)^T \theta)} \stackrel{\text{Defn. of } \sigma}{=} \frac{\exp(\phi(x_i)^T \theta)}{1 + \exp(\phi(x_i)^T \theta)} \stackrel{\text{Alg.}}{=} \frac{1}{1 + \exp(\phi(x_i)^T \theta)}$$

$\sigma(\phi(x)^T \theta) = \frac{1}{1 + \exp(-\phi(x)^T \theta)}$

$= \sigma(-\phi(x_i)^T \theta)$

Simplified Loss Minimization Problem

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta)))$$

The Loss for Logistic Regression

➤ Average **cross entropy** (simplified):

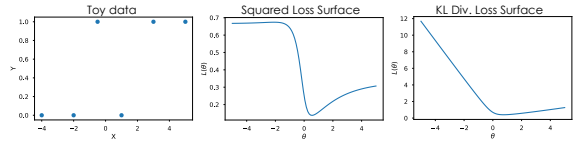
$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta)))$$

- Equivalent to (derived from) **minimizing the KL divergence**
- Also equivalent to **maximizing the log-likelihood of the data ...** (not covered in DS100 this semester)

Is this loss function reasonable?

Convexity Using Pictures

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta)))$$



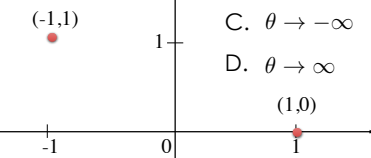
What is the value of θ ?

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta)))$$

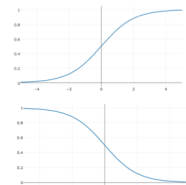
<http://bit.ly/ds100-sp18-cls>

Assume: $\phi(x) = x$

The Data



- A. $\theta = -1$
- B. $\theta = 1$
- C. $\theta \rightarrow -\infty$
- D. $\theta \rightarrow \infty$



What is the value of θ ?

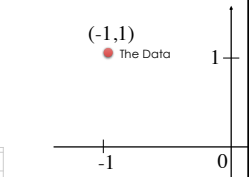
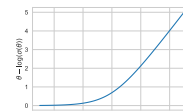
$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta)))$$

Assume: $\phi(x) = x$

For the point $(-1, 1)$: $y_i \phi(x_i)^T = -1$
 $-\phi(x_i)^T = 1$

Objective:

$$\theta - \log(\sigma(\theta))$$



$\theta \rightarrow -\infty$

What is the value of θ ?

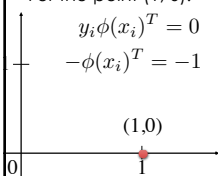
$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta)))$$

Assume: $\phi(x) = x$

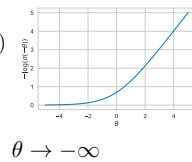
For the point $(1, 0)$:

$$y_i \phi(x_i)^T = 0 \rightarrow 0 - \log(\sigma(-\theta))$$

$$-\phi(x_i)^T = -1$$



For the point $(-1, 1)$: $\theta - \log(\sigma(\theta))$
 $\theta \rightarrow -\infty$



What is the value of θ ?

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta)))$$

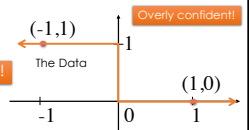
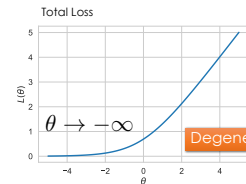
Assume: $\phi(x) = x$

For the point $(-1, 1)$: $\theta - \log(\sigma(\theta))$

$\theta \rightarrow -\infty$

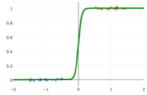
For the point $(1, 0)$: $0 - \log(\sigma(-\theta))$

$\theta \rightarrow -\infty$



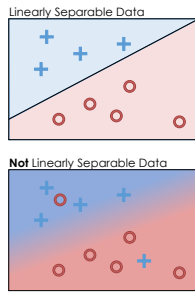
Linearly Separable Data

- > A classification dataset is said to be linearly separable if there exists a hyperplane that separates the two classes.
- > If data is linearly separable, logistic regression requires regularization



Weights go to infinity!

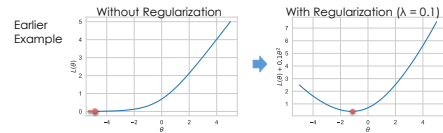
Solution?



Adding Regularization to Logistic Regression

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta))) + \lambda \sum_{j=1}^d \theta_j^2$$

- > Prevents weights from diverging on linearly separable data



Minimize the Loss

Logistic Loss Function

- > Average KL divergence (simplified)

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta)))$$

- > Take Derivative:

$$\begin{aligned} \nabla_{\theta} \mathbf{L}(\theta) &= -\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} y_i \phi(x_i)^T \theta + \nabla_{\theta} \log(\sigma(-\phi(x_i)^T \theta)) \\ &= -\frac{1}{n} \sum_{i=1}^n y_i \phi(x_i) + \nabla_{\theta} \log(\sigma(-\phi(x_i)^T \theta)) \end{aligned}$$

- > Average KL divergence (simplified)

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta)))$$

- > Take Derivative:

$$\begin{aligned} \nabla_{\theta} \mathbf{L}(\theta) &= -\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} y_i \phi(x_i)^T \theta + \nabla_{\theta} \log(\sigma(-\phi(x_i)^T \theta)) \\ &= -\frac{1}{n} \sum_{i=1}^n y_i \phi(x_i) + \nabla_{\theta} \log(\sigma(-\phi(x_i)^T \theta)) \\ &= -\frac{1}{n} \sum_{i=1}^n y_i \phi(x_i) + \frac{1}{\sigma(-\phi(x_i)^T \theta)} \nabla_{\theta} \sigma(-\phi(x_i)^T \theta) \end{aligned}$$

- > Take Derivative:

$$\nabla_{\theta} \mathbf{L}(\theta) = -\frac{1}{n} \sum_{i=1}^n y_i \phi(x_i) + \frac{1}{\sigma(-\phi(x_i)^T \theta)} \nabla_{\theta} \sigma(-\phi(x_i)^T \theta)$$

Useful Identity

$$\frac{\partial}{\partial t} \sigma(t) = \frac{\partial}{\partial t} \frac{1}{1 + e^{-t}} \stackrel{\text{Chain Rule}}{=} \frac{-1}{(1 + e^{-t})^2} \frac{\partial}{\partial t} (1 + e^{-t})$$

$$\stackrel{\text{Chain Rule}}{=} \frac{e^{-t}}{(1 + e^{-t})^2} \stackrel{\text{Alg.}}{=} \left(\frac{1}{1 + e^{-t}} \right) \left(\frac{e^{-t}}{1 + e^{-t}} \right)$$

$$\stackrel{\text{Alg.}}{=} \left(\frac{1}{1 + e^{-t}} \right) \left(\frac{1}{e^t + 1} \right) \stackrel{\text{Defn. of } \sigma}{=} \sigma(t) \sigma(-t)$$

➤ Take Derivative:

$$\nabla_{\theta} \mathbf{L}(\theta) = -\frac{1}{n} \sum_{i=1}^n y_i \phi(x_i) + \frac{1}{\sigma(-\phi(x_i)^T \theta)} \nabla_{\theta} \sigma(-\phi(x_i)^T \theta)$$

Useful Identity $\frac{\partial}{\partial t} \sigma(t) = \sigma(t)\sigma(-t)$

$$= -\frac{1}{n} \sum_{i=1}^n y_i \phi(x_i) + \frac{\sigma(-\phi(x_i)^T \theta)}{\sigma(-\phi(x_i)^T \theta)} \sigma(\phi(x_i)^T \theta) \nabla_{\theta} (-\phi(x_i)^T \theta)$$

$$= -\frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\phi(x_i)^T \theta)) \phi(x_i)$$

Logistic Loss Function

➤ Average KL divergence (simplified)

$$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^n (y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta)))$$

➤ Take Derivative:

$$\nabla_{\theta} \mathbf{L}(\theta) = -\frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\phi(x_i)^T \theta)) \phi(x_i)$$

➤ Set derivative = 0 and solve for θ

- No general analytic solution
- Solved using numeric methods

The Gradient Descent Algorithm

$\theta^{(0)} \leftarrow$ initial vector (random, zeros ...)

For τ from 0 to convergence:

$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \rho(\tau) \left(\nabla_{\theta} \mathbf{L}(\theta) \Big|_{\theta=\theta^{(\tau)}} \right)$$

➤ $\rho(\tau)$ is the step size (learning rate)

- typically $1/\tau$
- Converges when gradient is ≈ 0 (or we run out of patience)

Gradient Descent for Logistic Regression

Logistic Regression

$\theta^{(0)} \leftarrow$ initial vector (random, zeros ...)

For τ from 0 to convergence:

$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \rho(\tau) \left(\frac{1}{n} \sum_{i=1}^n (\sigma(\phi(x_i)^T \theta^{(\tau)}) - y_i) \phi(x_i) \right)$$

➤ $\rho(\tau)$ is the step size (learning rate)

- typically $1/\tau$
- Converges when gradient is ≈ 0 (or we run out of patience)

Stochastic Gradient Descent

➤ For many learning problems the gradient is a sum:

$$\nabla_{\theta} \mathbf{L}(\theta) = \frac{1}{n} \sum_{i=1}^n (\sigma(\phi(x_i)^T \theta) - y_i) \phi(x_i)$$

➤ For large n this can be costly

➤ What if we approximated the gradient by looking at a few random points:

$$\nabla_{\theta} \mathbf{L}(\theta) \approx \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} (\sigma(\phi(x_i)^T \theta) - y_i) \phi(x_i)$$

➤ What if we approximated the gradient by looking at a few random points:

$$\nabla_{\theta} \mathbf{L}(\theta) \approx \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} (\sigma(\phi(x_i)^T \theta) - y_i) \phi(x_i)$$

Batch Size

Random sample of records

➤ This is a reasonable estimator for the gradient

- Unbiased ...
- Often batch size is one! (why is this helpful)
 - Fast to compute!
- A key ingredient in the recent success of deep learning

Stochastic Gradient Descent

$\theta^{(0)} \leftarrow$ initial vector (random, zeros ...)

For τ from 0 to convergence:

$\mathcal{B} \sim$ Random subset of indices

$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \rho(\tau) \left(\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla_{\theta} \mathbf{L}_i(\theta) \Big|_{\theta=\theta^{(\tau)}} \right)$$

Decomposable Loss $\mathbf{L}(\theta) = \sum_{i=1}^n \mathbf{L}_i(\theta) = \sum_{i=1}^n \mathbf{L}(\theta, x_i, y_i)$

Loss can be written as a sum of the loss on each record.

$\theta^{(0)} \leftarrow$ initial vector (random, zeros ...)

For τ from 0 to convergence:

$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \rho(\tau) \left(\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \mathbf{L}_i(\theta) \Big|_{\theta=\theta^{(\tau)}} \right)$$

Gradient Descent

Assuming Decomposable Loss Functions

$\theta^{(0)} \leftarrow$ initial vector (random, zeros ...)

For τ from 0 to convergence:

$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \rho(\tau) \left(\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla_{\theta} \mathbf{L}_i(\theta) \Big|_{\theta=\theta^{(\tau)}} \right)$$

Stochastic Gradient Descent

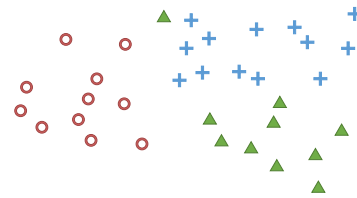
Very Similar Algorithms

$\mathcal{B} \sim$ Random subset of indices

Python Demo!

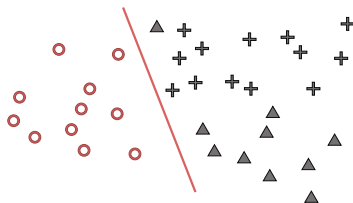
Multiclass (more than 2) Classification

➤ **One-vs-rest** train separate binary classifiers for each class



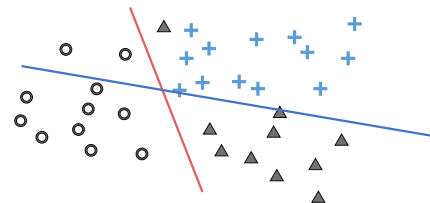
Multiclass (more than 2) Classification

➤ **One-vs-rest** train separate binary classifiers for each class



Multiclass (more than 2) Classification

➤ **One-vs-rest** train separate binary classifiers for each class



Multiclass (more than 2) Classification

- **One-vs-rest** train separate binary classifiers for each class

Multiclass (more than 2) Classification

- **One-vs-rest** train separate binary classifiers for each class

Multiclass (more than 2) Classification

- **One-vs-rest** train separate binary classifiers for each class
 - Class with highest confidence wins
 - Need to address class imbalance issue
- **Soft-Max** multiclass classification

- **Soft-Max** multiclass classification

$$\mathbf{P}(Y = j | x) = \frac{\exp(x^T \theta^{(j)})}{\sum_{m=1}^k \exp(x^T \theta^{(m)})}$$

- Separate $\theta^{(j)} \in \mathbb{R}^p$ for each class
- Trained using gradient descent methods
- Over parameterized. Why?
 - k sets of parameters one for each class
 - Only need k-1 parameters

$$\mathbf{P}(y = k | x) = 1 - \sum_{j=1}^K \mathbf{P}(y = j | x)$$

- Often use k parameters + regularization to address "redundancy".

Python Demo!

Deep Learning Overview

Bonus Material

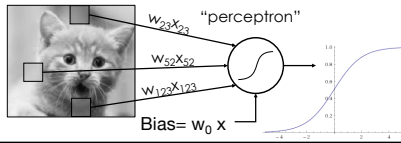
Borrowed heavily from excellent talks by:

- Adam Coates: http://ai.stanford.edu/~acoates/coates_dltutorial_2013.pptx
- Fei-Fei Li and Andrej Karpathy: <http://cs231n.stanford.edu/syllabus.html>

Logistic Regression as a "Neuron"

Consider the simple function family: $\sigma(u) = \frac{1}{1 + \exp(-u)}$

$$f_w(x) = \sigma(w^T x) = \sigma\left(\sum_{j=1}^d w_j x_j\right) = P(y = 1 | x)$$



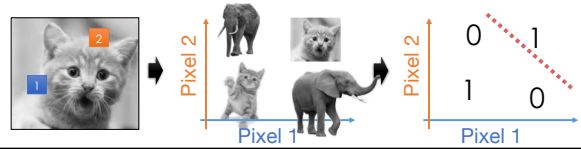
Neuron "fires" if weighted sum of input is greater than zero.

Logistic Regression: Strengths and Limitations

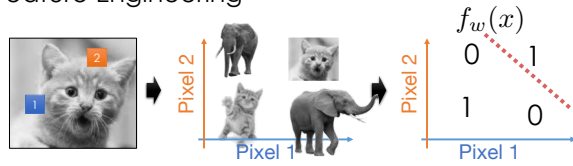
- Widely used machine learning technique
- convex → efficient to learn
- easy to interpret model weights
- works well given good features



- Limitations:
- Restricted to linear relationships → sensitive to choice of features



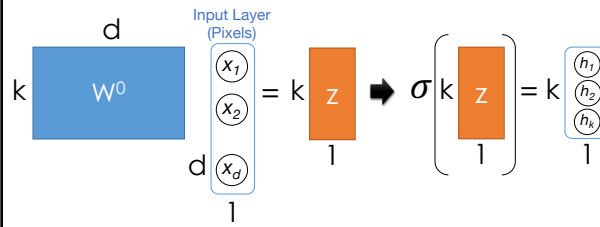
Feature Engineering



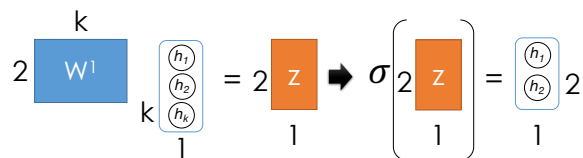
Rather than use raw pixels build/train feature functions:



Composition Linear Models and Nonlinearities



Composition Linear Models and Nonlinearities

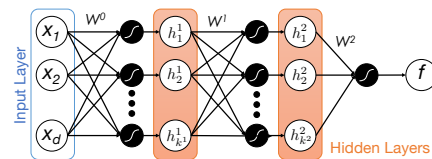


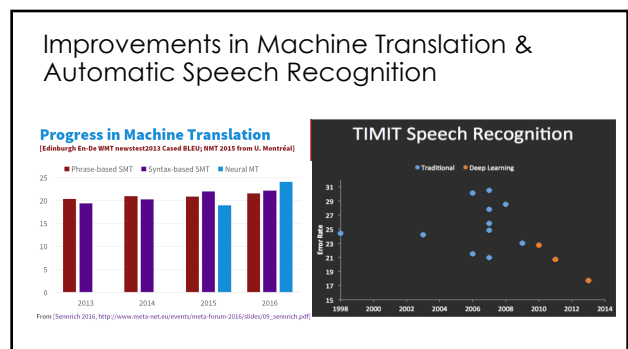
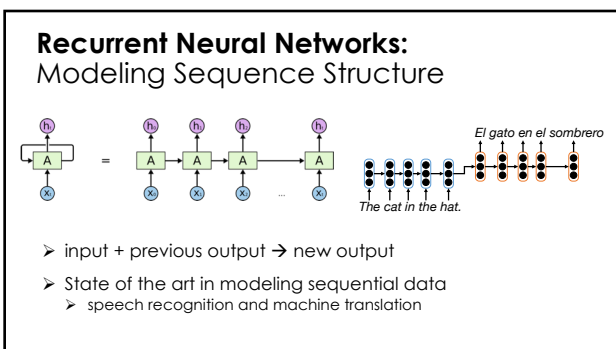
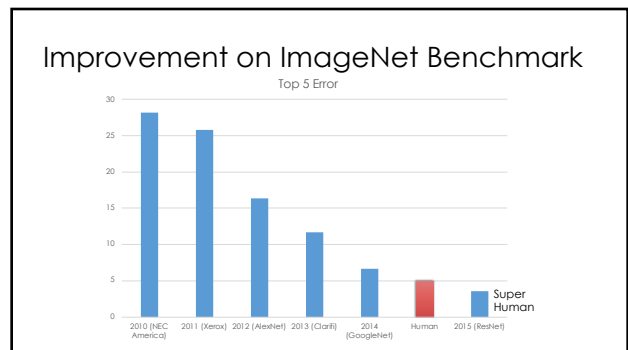
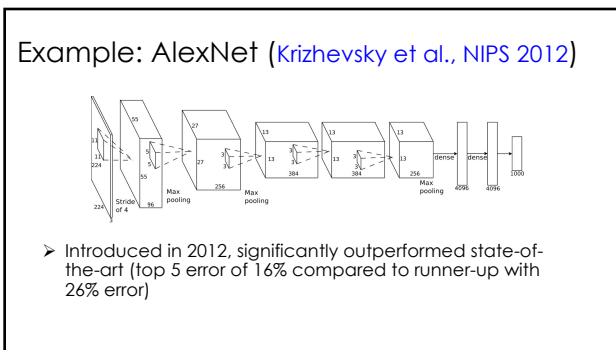
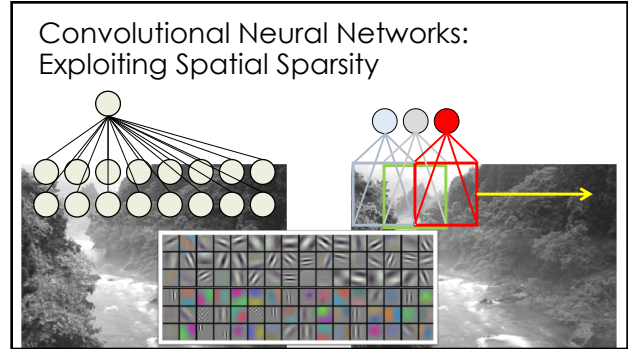
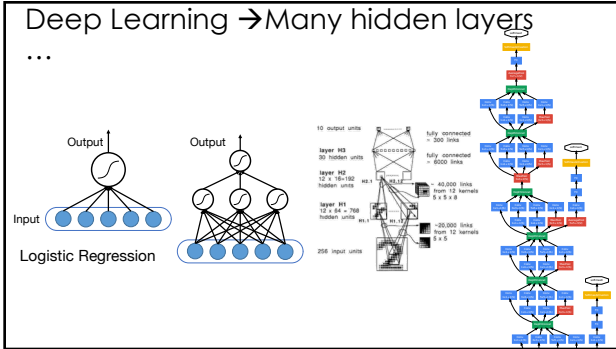
Neural Networks

Composing "perceptrons"

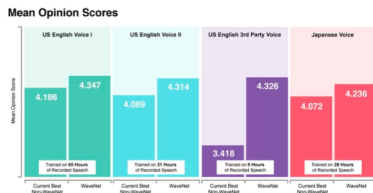
$$x \rightarrow \sigma(W^0 x) \rightarrow h^1 \rightarrow \sigma(W^1 h^1) \rightarrow h^2 \rightarrow \sigma(W^2 h^2) \rightarrow f$$

$$y = f_{W^0, W^1, W^2}(x) = \sigma(W^2 \sigma(W^1 \sigma(W^0 x)))$$





State of the art in Text to Speech (TTS)



Interested in Deep Learning?

- RISE Lab Deep Learning Overview:
 - https://ucbrise.github.io/cs294-rise-fa16/deep_learning.html
- [TensorFlow Python Tutorial](#)
- Stanford CS231 Labs
 - <http://cs231n.github.io/linear-classify/>
 - <http://cs231n.github.io/optimization-1/>
 - <http://cs231n.github.io/optimization-2/>