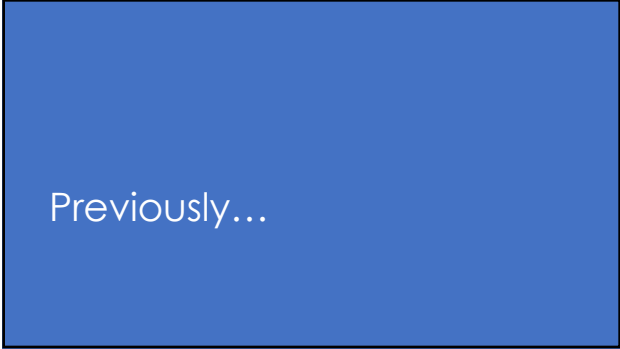
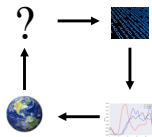


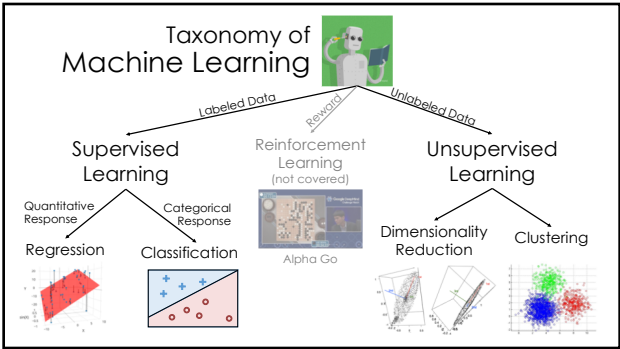
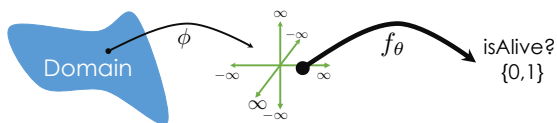
Classification & Logistic Regression & maybe deep learning

Slides by:

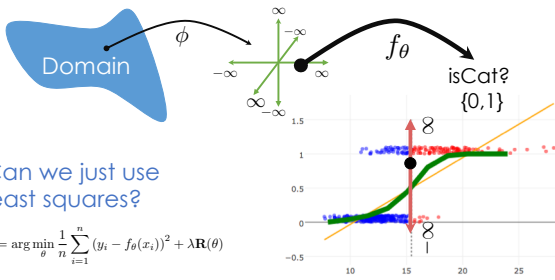
Joseph E. Gonzalez jgonzal@cs.berkeley.edu



Classification



Classification

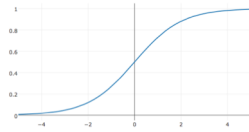


Logistic Regression

➤ Model the probability of a particular label:

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

$$\sigma(t) = \frac{1}{1 + \exp(-t)}$$



Motivation for the Logistic Model

➤ Model the "log odds" as a linear model

$$\underbrace{\phi(x_i)^T \theta}_{\text{Linear Model}} = \log \left(\underbrace{\frac{\hat{P}_\theta(y_i = 1 | x_i)}{\hat{P}_\theta(y_i = 0 | x_i)}}_{\text{Log odds}} \right)$$

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

$$\phi(x_i)^T \theta > 0 \quad \overset{\text{exp}(t) > 1}{\Rightarrow} \quad \hat{P}_\theta(y_i = 1 | x_i) > \hat{P}_\theta(y_i = 0 | x_i)$$

$$\phi(x_i)^T \theta < 0 \quad \overset{\text{exp}(-t) < 1}{\Rightarrow} \quad \hat{P}_\theta(y_i = 1 | x_i) < \hat{P}_\theta(y_i = 0 | x_i)$$

for any positive t

Motivation for the Logistic Model

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

Taking the exponent of both sides

$$\exp(\phi(x_i)^T \theta) = \frac{\hat{P}_\theta(y_i = 1 | x_i)}{1 - \hat{P}_\theta(y_i = 1 | x_i)}$$

$$\exp(\phi(x_i)^T \theta) = \frac{\hat{P}_\theta(y_i = 1 | x_i)}{1 - \hat{P}_\theta(y_i = 1 | x_i)}$$

Algebra

$$\exp(\phi(x_i)^T \theta) (1 - \hat{P}_\theta(y_i = 1 | x_i)) = \hat{P}_\theta(y_i = 1 | x_i)$$

Expanding terms

$$\exp(\phi(x_i)^T \theta) - \exp(\phi(x_i)^T \theta) \hat{P}_\theta(y_i = 1 | x_i) = \hat{P}_\theta(y_i = 1 | x_i)$$

Collect terms on the other side ...

$$\exp(\phi(x_i)^T \theta) = \hat{P}_\theta(y_i = 1 | x_i) (1 + \exp(\phi(x_i)^T \theta))$$

Solving for P(y=1|x)

$$\hat{P}_\theta(y_i = 1 | x_i) = \frac{\exp(\phi(x_i)^T \theta)}{1 + \exp(\phi(x_i)^T \theta)}$$

Solving for P(y=1|x)

$$\hat{P}_\theta(y_i = 1 | x_i) = \frac{\exp(\phi(x_i)^T \theta)}{1 + \exp(\phi(x_i)^T \theta)}$$

Dividing numerator and denominator by exp(phi(x)^T theta)

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

$$= \sigma(\phi(x)^T \theta)$$

Where

$$\sigma(t) = \frac{1}{1 + \exp(-t)}$$

The Logistic Regression Model

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

How do we fit the model to the data?

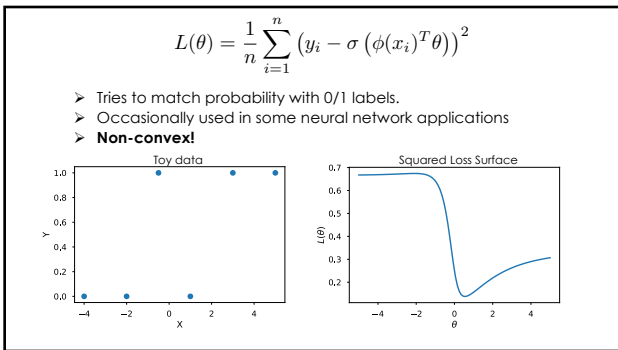
Defining the Loss

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!**




Defining the Cross Entropy Loss

Loss Function

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

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

- Example: (cute or not)?

x = 

y = 1 "cute"

	Cute	Not Cute
Observed Probability	$\mathbf{P}(y = 1 x) = 1.0$	$\mathbf{P}(y = 0 x) = 0.0$
Predicted Probability	$\hat{\mathbf{P}}_{\theta}(y = 1 x) = 0.8$	$\hat{\mathbf{P}}_{\theta}(y = 0 x) = 0.2$

Loss Function

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

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

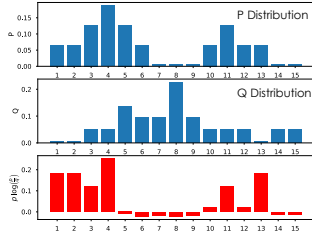
- Kullback-Leibler (KL) Divergence provides a measure of difference between two distributions:
- Difference between two discrete distributions P and Q

$$\mathbf{D}(P||Q) = \sum_{k=0}^{K-1} P(k) \log \left(\frac{P(k)}{Q(k)} \right)$$

Kullback–Leibler (KL) Divergence

$$D(P||Q) = \sum_{k=0}^{K-1} P(k) \log \left(\frac{P(k)}{Q(k)} \right)$$

- The average log difference between P and Q weighted by P
- Does not penalize mismatch for rare events with respect to P
- Note that it is not symmetric
 $D(P||Q) \neq D(Q||P)$



Loss Function

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

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

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

$$D_{KL}(P||\hat{P}_\theta) = \sum_{k=0}^{K-1} P(y = k | x) \log \left(\frac{P(y = k | x)}{\hat{P}_\theta(y = k | x)} \right)$$

- Average KL Divergence for all the data:

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

$$D_{KL}(P||\hat{P}_\theta) = \sum_{k=0}^{K-1} P(y = k | x) \log \left(\frac{P(y = k | x)}{\hat{P}_\theta(y = k | x)} \right)$$

- Average KL Divergence for all the data:

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

Doesn't depend on θ

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

- Average **cross entropy loss**

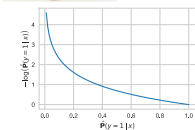
$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \sum_{k=0}^{K-1} - P(y_i = k | x_i) \log(\hat{P}_\theta(y_i = k | x_i))$$

- Average **cross entropy loss**

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \sum_{k=0}^{K-1} - P(y_i = k | x_i) \log(\hat{P}_\theta(y_i = k | x_i))$$

Cute Cat Example

x =  , y = 1 "cute"



	Cute	Not Cute
Observed Probability	$P(y = 1 x) = 1.0$	$P(y = 0 x) = 0.0$
Predicted Probability	$\hat{P}_\theta(y = 1 x) = 0.8$	$\hat{P}_\theta(y = 0 x) = 0.2$
Cross Ent. - $P \log \hat{P}_\theta$	$-1.0 \log(0.8) \approx 0.22$	$-0.0 \log(0.2) = 0.0$

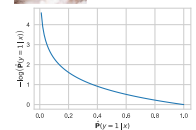
Also called the log loss because it is the log of the predicted probability for the true class

- Average **cross entropy loss**

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \sum_{k=0}^{K-1} - P(y_i = k | x_i) \log(\hat{P}_\theta(y_i = k | x_i))$$

Cute Cat Example

x =  , y = 0 "not cute"



	Cute	Not Cute
Observed Probability	$P(y = 1 x) = 0.0$	$P(y = 0 x) = 1.0$
Predicted Probability	$\hat{P}_\theta(y = 1 x) = 0.7$	$\hat{P}_\theta(y = 0 x) = 0.3$
Cross Ent. - $P \log \hat{P}_\theta$	$-0.0 \log(0.7) = 0.0$	$-1.0 \log(0.3) \approx 1.20$

Also called the log loss because it is the log of the predicted probability for the true class

➤ Average **cross entropy loss**

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \sum_{k=0}^{K-1} -\mathbf{P}(y_i = k | x_i) \log(\hat{\mathbf{P}}_{\theta}(y_i = k | x_i))$$

➤ Computing the more general version for (x_i, y_i)

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

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

➤ Average **cross entropy loss**

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \sum_{k=0}^{K-1} -\mathbf{P}(y_i = k | x_i) \log(\hat{\mathbf{P}}_{\theta}(y_i = k | x_i))$$

➤ Computing the more general version for (x_i, y_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**

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \sum_{k=0}^{K-1} -\mathbf{P}(y_i = k | x_i) \log(\hat{\mathbf{P}}_{\theta}(y_i = k | x_i))$$

Rewriting on one line:

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

After much algebra (see last lecture) we obtain:

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

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 Data100 this semester)

Is this loss function reasonable?

Convexity Using Pictures

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

Toy data

Squared Loss Surface

KL Div. Loss Surface

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: $\phi(x) = x$

The Data

$(-1, 1)$

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

For the point (1, 0): $y_i \phi(x_i)^T = 0$ $\rightarrow 0 - \log(\sigma(-\theta))$
 $-\phi(x_i)^T = -1$

$\theta \rightarrow -\infty$

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): $0 - \log(\sigma(-\theta))$
 $\theta \rightarrow -\infty$

Total Loss

$\theta \rightarrow -\infty$

Degenerate Solution!

Overly confident!

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

Earlier Example

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:

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

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

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

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

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:

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

- $\rho(\theta)$ 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.

Gradient Descent

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

For τ from 0 to convergence:

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

Stochastic Gradient Descent

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

For τ from 0 to convergence:

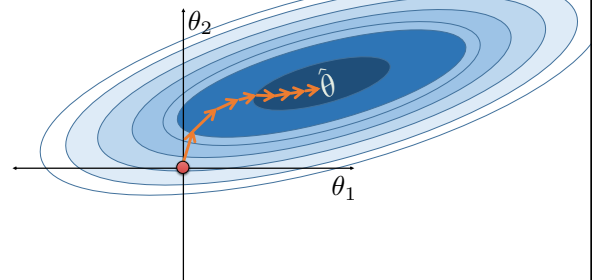
$\mathcal{B} \sim$ Random subset of indices

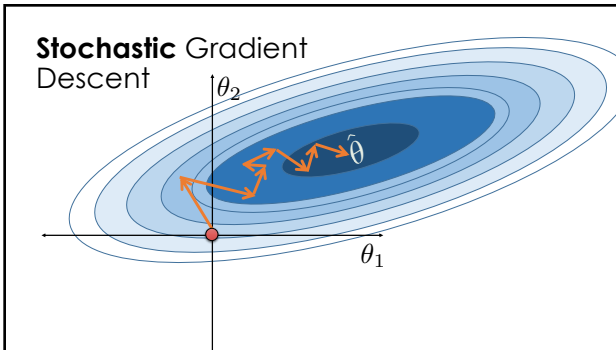
$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \rho(\tau) \left(\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \nabla_{\theta} \mathbf{L}_i(\theta) \Big|_{\theta=\theta^{(\tau)}} \right)$$

Very Similar Algorithms

Assuming Decomposable Loss Functions

Gradient Descent





Logistic Regression in Scikit Learn

```

from sklearn.linear_model import LogisticRegression
# By default SK learn adds regularization
# C = 1/lambda the inverse regularization parameter.
model = LogisticRegression(C=100.00)
# Train the model
model.fit(df[['feat1', 'feat2']], df['label'])
# Make Predictions
test_df['label'] = model.predict(test_df[['feat1', 'feat2']])
test_df['P(Y|X)'] = model.predict_proba(test_df[['feat1', 'feat2']])

```

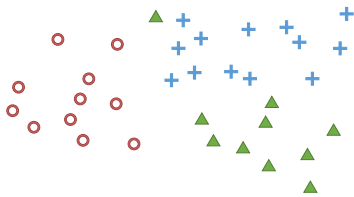
Python Demo!

Attendance Quiz <http://bit.ly/ds100-sp18-sgd>

1. Is the gradient of a **simple random sample** of the data an **unbiased estimate** of the gradient of the entire dataset? (T/F)
2. By decreasing the batch size we:
 1. **Increase** the **variance** (T/F)
 2. **Decrease** the **bias** (T/F)
 3. **Reduce** the **computational cost** of each iteration (T/F)

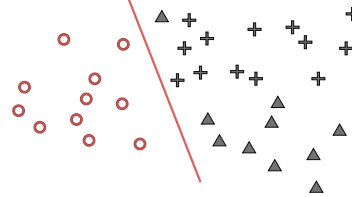
Multiclass (more than 2) Classification

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



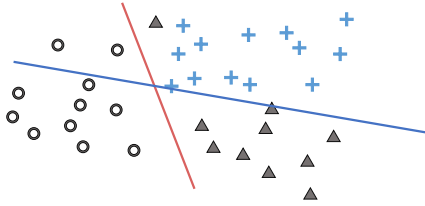
Multiclass (more than 2) Classification

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



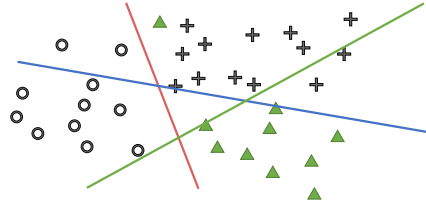
Multiclass (more than 2) Classification

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



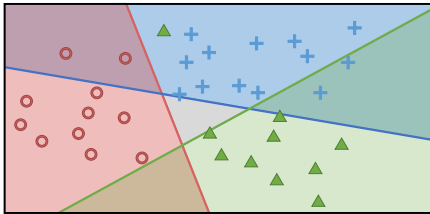
Multiclass (more than 2) Classification

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



Multiclass (more than 2) Classification

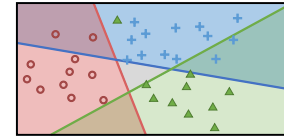
- **One-vs-rest** train separate binary classifiers for each class



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

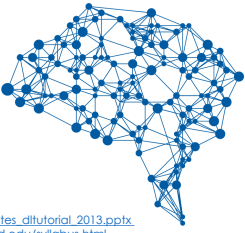
$$\mathbf{P}(y = k | x) = 1 - \sum_{j=1}^K \mathbf{P}(y = j | x)$$

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

Python Demo!

Deep Learning Overview

Bonus Material

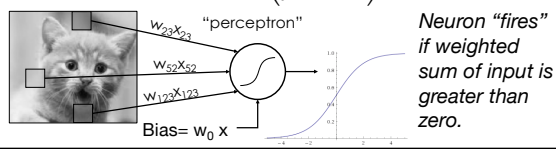


Borrowed from excellent talks by:

- 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"

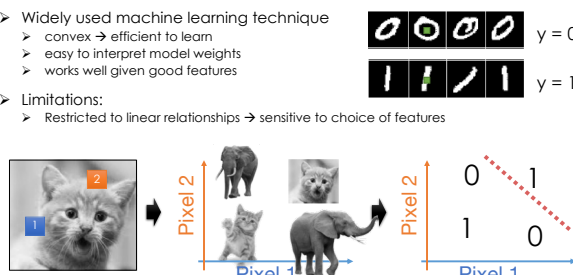
➤ Consider the simple function family: $\sigma(u) = \frac{1}{1 + \exp(-u)}$

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


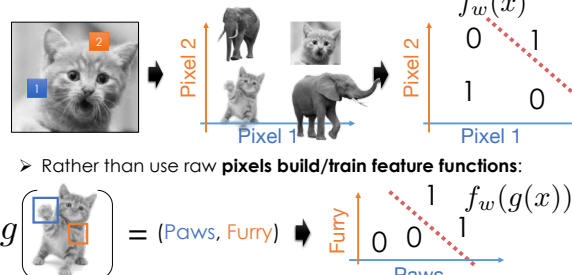
Neuron "fires" if weighted sum of input is greater than zero.

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



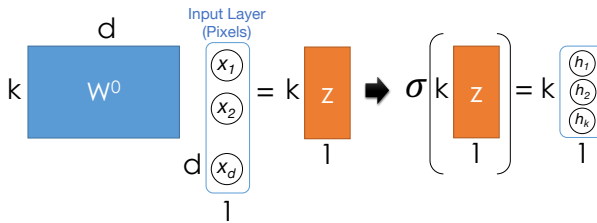
Feature Engineering



➤ Rather than use raw pixels build/train feature functions:

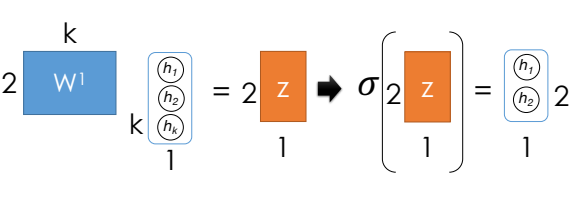
$$g\left(\begin{matrix} \text{Kitten} \\ \text{Kitten} \end{matrix}\right) = (\text{Paws}, \text{Furry}) \Rightarrow \begin{matrix} \text{Furry} \\ \text{Paws} \end{matrix}$$

Composition Linear Models and Nonlinearities



$$\begin{matrix} d \\ k \end{matrix} W^0 \begin{matrix} \text{Input Layer (Pixels)} \\ x_1 \\ x_2 \\ \vdots \\ x_d \\ 1 \end{matrix} = \begin{matrix} k \\ 1 \end{matrix} z \Rightarrow \sigma \left[\begin{matrix} k \\ 1 \end{matrix} z \right] = \begin{matrix} \hat{h}_1 \\ \hat{h}_2 \\ \vdots \\ \hat{h}_k \\ 1 \end{matrix}$$

Composition Linear Models and Nonlinearities



$$\begin{matrix} k \\ 2 \end{matrix} W^1 \begin{matrix} \hat{h}_1 \\ \hat{h}_2 \\ \vdots \\ \hat{h}_k \\ 1 \end{matrix} = \begin{matrix} 2 \\ 1 \end{matrix} z \Rightarrow \sigma \left[\begin{matrix} 2 \\ 1 \end{matrix} z \right] = \begin{matrix} \hat{h}_1 \\ \hat{h}_2 \\ 1 \end{matrix}$$

Neural Networks

- Composing "perceptrons"

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

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

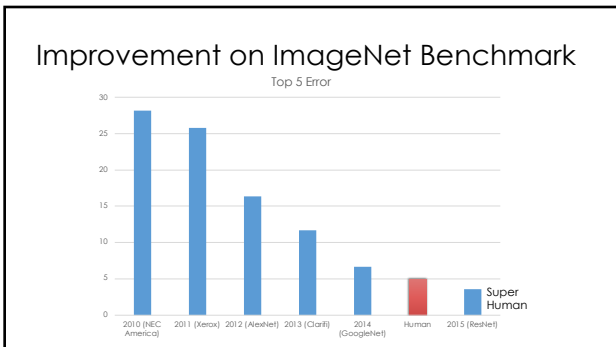
Deep Learning → Many hidden layers

...

Convolutional Neural Networks: Exploiting Spatial Sparsity

Example: AlexNet (Krizhevsky et al., NIPS 2012)

- Introduced in 2012, significantly outperformed state-of-the-art (top 5 error of 16% compared to runner-up with 26% error)



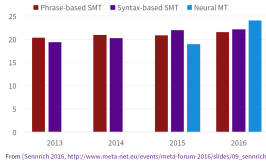
Recurrent Neural Networks: Modeling Sequence Structure

- input + previous output → new output
- State of the art in modeling sequential data
 - speech recognition and machine translation

Improvements in Machine Translation & Automatic Speech Recognition

Progress in Machine Translation

[Edinburgh En-De WMT newstest2013 Case4 BLEU; NMT 2015 from U. Montréal]



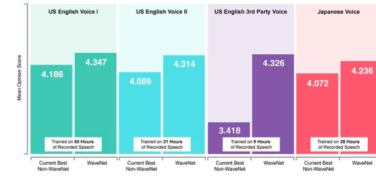
From [Sennrich 2016, http://www.meta-net.eu/events/meta_forum_2016/slides/01_sennrich.pdf]

TIMIT Speech Recognition



State of the art in Text to Speech (TTS)

Mean Opinion Scores



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/>