

Softmax (Logistic) Regression

Linear regression: Predict a continuous value.

Classification: Predict one of k classes.

- classes might not be ordered

~~"cat" > "dog"~~

- might want to predict the probability for each possible class.

ex 93% cat
 7% dog

ex Predict sell above, below or not sell
from size, age, sq ft, price

"one-hot" vectors:

$[1, 0, 0]$	$[0, 1, 0]$	$[0, 0, 1]$
below	above	not sell
1	2	3

Model predict $k=3$ scores.

Larger score: "I think it's this class"

$$\underset{\substack{\text{score} \\ \downarrow \\ \text{class 1}}}{O_1} = \underset{\substack{\text{input} \\ \downarrow \\ \text{output size}}}{W_{1,1}} X_1 + W_{1,2} X_2 + \dots + b_1$$

age

$$O = Wx + b \quad W \in \mathbb{R}^{3 \times 4}, \quad b \in \mathbb{R}^3, \quad x \in \mathbb{R}^4$$
$$O \in \mathbb{R}^3$$

have unnormalized scores (any real number)
want probabilities:

1. Sum to 1

2. Lie in $[0, 1]$

$$\hat{y} = \text{softmax}(o) \quad o \in \mathbb{R}^k \quad \hat{y} \in \mathbb{R}^k$$

$$\text{softmax}(o)_j = \frac{\exp(o_j)}{\sum_k \exp(o_k)}$$

$$1. \sum_j \text{softmax}(o)_j = \frac{\sum_j \exp(o_j)}{\sum_k \exp(o_k)} = 1$$

$$2. \exp(\dots) \geq 0$$

argmax gives us the highest-probability class

$$\text{argmax}(o) = \text{argmax}(\text{softmax}(o))$$

softmax regression

$$p(\underset{\substack{\text{class} \\ \text{label}}}{y} | \underset{\text{input}}{x}) = \text{softmax}(\underset{\text{params}}{Wx + b})$$

"linear" because the ^{→0}scores are a linear function of the input

Want to find parameters that
maximize $p(y|x)$ over training dataset.

↳ ground-truth label

$$\prod_i p(y^{(i)} | x^{(i)}) \quad (x^{(i)}, y^{(i)} : i^{\text{th}} \text{ training example})$$

$$\sum_i \log p(y^{(i)} | x^{(i)}) \quad (\text{because log monotone})$$

$$- \sum_i \log p(y^{(i)} | x^{(i)}) \quad (\text{because we minimize})$$

$$- \sum_i \sum_j y_j^{(i)} \log \overbrace{\text{softmax}(W x^{(i)} + b)}^{\text{model}}_j$$

"cross-entropy"

measures how "wrong" a probability distribution is

$$-\sum_i \sum_j y_j^{(i)} \log \hat{y}_j^{(i)}$$

$$\hat{y}_j^{(i)} = \text{softmax}(W x^{(i)} + b)_j$$

for one particular example:

$$L(y, \hat{y}) = -\sum_j y_j \log \frac{\exp(o_j)}{\sum_k \exp(o_k)}$$

$$= \sum_j y_j \log \sum_k \exp(o_k) - \sum_j y_j o_j$$

$$= \log \sum_k \exp(o_k) - \sum_j y_j o_j$$

$$\frac{dL(\gamma, \hat{y})}{d o_j} \rightarrow \text{only new term}$$

$$\begin{aligned} d o_j [L(\gamma, \hat{y})] &= d o_j [\log \sum_k \exp(o_k) - \sum_j y_j o_j] \\ &= \frac{d o_j \sum_k \exp(o_k)}{\sum_k \exp(o_k)} = y_j \\ &= \frac{\exp(o_j)}{\sum_k \exp(o_k)} - y_j \\ &= \text{softmax}(o)_j - y_j \end{aligned}$$

Gradient Descent

We want to minimize $f(x)$

$$f(x+\epsilon) = \sum_{n=0}^{\infty} \frac{\epsilon^n f^{(n)}(x)}{n!}$$

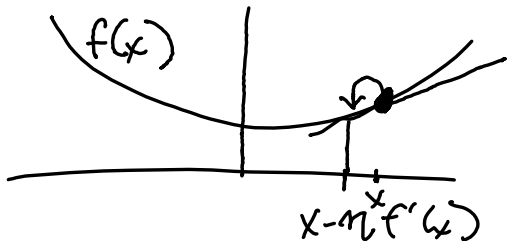
$$f(\underset{\substack{\uparrow \\ \text{small}}}{x+\epsilon}) = f(x) + \epsilon f'(x) + O(\epsilon^2)$$

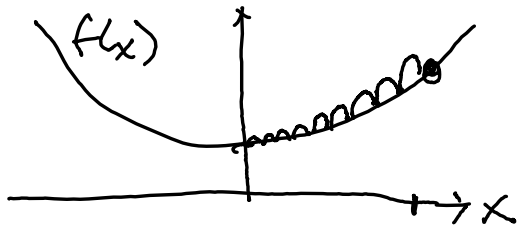
set $\epsilon = -\eta f'(x)$ η small, positive

$$f(x - \eta f'(x)) = f(x) - \underbrace{\eta (f'(x))^2}_{\text{non-negative}} + O(\eta^2 (f'(x))^2)$$

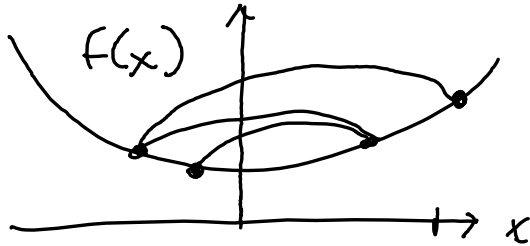
$$f(x - \eta f'(x)) \lesssim f(x)$$

Therefore: $x \leftarrow x - \eta f'(x)$

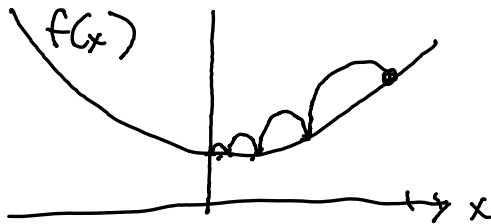




small η



big η



"just right" η

In the multivariate

$$\nabla_x f(x) = \underset{\text{vector}}{\left[\frac{\partial f(x)}{\partial x_1}, \frac{\partial f(x)}{\partial x_2}, \dots \right]^T}$$

$$f(x + \varepsilon) = f(x) + \varepsilon^T \nabla f(x) + O(\|\varepsilon\|^2)$$

$$x \leftarrow x - \underset{\text{scalar}}{\eta} \nabla f(x)$$

In ML, we have $L(y^{(i)}, x^{(i)}, \theta) = L_i(\theta)$
we want to minimize L by changing θ

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n L_i(\theta)$$

$$\nabla_{\theta} L(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L_i(\theta)$$

"Batch" gradient descent

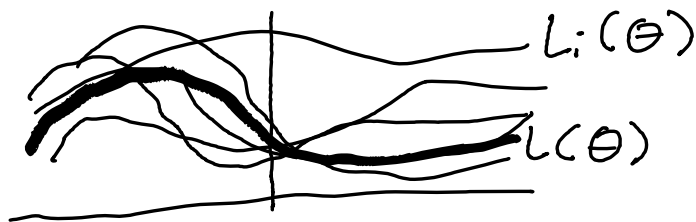
Stochastic gradient descent:

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} L_i(\Theta)$$

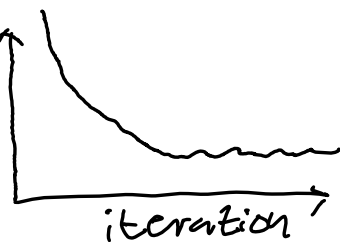
$i \sim \text{uniform categorical}(n)$

$$\mathbb{E}_i[\nabla_{\Theta} L_i(\Theta)] = \frac{1}{n} \sum_{i=1}^n \nabla_{\Theta} L_i(\Theta)$$

We are optimizing a different function at each iteration.



Progress can stall as we approach minimum of $L(\Theta)$



One option: Use a "learning rate schedule"

→ different η at each iteration t

$$\eta(t) = \eta \quad (\text{constant})$$

$$\eta(t) = \eta_i \quad \text{if } t_i \leq t \leq t_{i+1} \quad (\text{piecewise})$$

$$\eta(t) = \eta_0 e^{-\alpha t} \quad (\text{exponential})$$

Another option: Use a minibatch.

Sample B examples randomly and use their average gradient.

$$\Theta \leftarrow \Theta - \eta \frac{1}{B} \sum_{i=1}^B \nabla_{\Theta} L_i(\Theta)$$

Note: SGD can be more efficient BGD

if $\nabla_{\Theta} L_i$ is generally aligned with $\nabla_{\Theta} L(\Theta)$

Note: SGD can be noisy.

Note: SGD is less parallelizable.

Why might it be a bad idea to use the same η for every parameter?

ex $\hat{y} = w^T x$

$$L(y, \hat{y}) = \frac{1}{2} \|y - \hat{y}\|^2$$

$$\frac{dL}{d\hat{y}} = \hat{y} - y$$

$$\frac{dL}{dw} = (\hat{y} - y) x$$

Imagine that $|x_i| \gg |x_j|$ e.g.

$x_i \sim \mathcal{N}(0, \sigma_i)$ $x_j \sim \mathcal{N}(0, \sigma_j)$ with $\sigma_i \gg \sigma_j$

Then we will often have $|\frac{\partial L}{\partial w_i}| \gg |\frac{\partial L}{\partial w_j}|$