

ECE C147 Midterm Cheat Sheet

$$\nabla_W \left(\frac{\alpha}{2} W^T W \right) = \alpha W$$

L2 Regularization

$$\tilde{J}(W; X, y) = J(W; X, y) + \frac{\alpha}{2} W^T W$$

$$W \leftarrow (I - \epsilon I) W - \epsilon \nabla_W J(W; X, y)$$

ML basics

K-nearest → too much memory, testing is expensive
→ distance measurements in high-D space do not reflect similarity

• calculate distance between $x^{(i)}$ & $x^{(j)}$ & c

• take the K nearest indices, classify via mode of $\{y^{(1)}, y^{(2)}, \dots, y^{(n)}\}$

Prob: Simple

Cons: Slow/scale w/ amount of training data, doesn't work if dimensionality is high, near neighbors are far away
L2 dist between rows: $\sqrt{(x_{1,1})^2 + \dots + (x_{1,n})^2} = \sqrt{x_1^2 + \dots + x_n^2} = \sqrt{2n \cdot \text{var}(x_i)}$

Broadcasting

class edges	one has	size 1
(2x3)	(2x1)	(3)
(3)	(3)	
ax3	3x3	

$$\mathbb{E}[\delta \delta^T] = 0^3, \text{ If } \mathbb{E}[\delta] = 0$$

Softmax

$$\text{softmax}_i(x) = \frac{e^{w_i^T x + b_i}}{\sum_j e^{w_j^T x + b_j}}, \quad \text{probability that } x \text{ belongs to class } i,$$

$$P(y^{(i)} = i \mid x^{(i)}, \theta) = \text{softmax}_i(x^{(i)})$$

$$\begin{aligned} \Theta^T = \arg\max_{\Theta} \sum_{i=1}^m \log \text{softmax}_i(x^{(i)}) \\ = \arg\max_{\Theta} \sum_{i=1}^m \log \left[\frac{e^{w_i^T x^{(i)} + b_i}}{\sum_j e^{w_j^T x^{(i)} + b_j}} \right] \\ = \arg\max_{\Theta} \sum_{i=1}^m \left[\log \left(\frac{e^{w_i^T x^{(i)} + b_i}}{\sum_j e^{w_j^T x^{(i)} + b_j}} \right) - \log \left(\sum_j e^{w_j^T x^{(i)} + b_j} \right) \right] \\ = \arg\max_{\Theta} \sum_{i=1}^m \left[\log \left(\frac{e^{w_i^T x^{(i)} + b_i}}{\sum_j e^{w_j^T x^{(i)} + b_j}} \right) - c_i \right] \end{aligned}$$

$L = \log \text{likelihood}$

$$\nabla_{\Theta} L = \nabla_{\Theta} \left(\frac{1}{m} \sum_{i=1}^m \left[\log \left(\frac{e^{w_i^T x^{(i)} + b_i}}{\sum_j e^{w_j^T x^{(i)} + b_j}} \right) - c_i \right] \right)$$

the class of the j^{th} sample

how many times does $y^{(i)} = i$?

Let c_i represent the # of samples that belong to class i .

$$\nabla_{w_i} (\log \left(\frac{e^{w_i^T x^{(i)} + b_i}}{\sum_j e^{w_j^T x^{(i)} + b_j}} \right)) = \frac{\sum_{k=1}^m \nabla_{w_k} (e^{w_k^T x^{(i)} + b_k})}{\sum_{k=1}^m e^{w_k^T x^{(i)} + b_k}} = \frac{e^{w_i^T x^{(i)} + b_i}}{\sum_{k=1}^m e^{w_k^T x^{(i)} + b_k}}$$

$$\therefore \frac{1}{m} \left[\sum_{j=1}^m \left(\frac{e^{w_j^T x^{(i)} + b_j}}{\sum_{k=1}^m e^{w_k^T x^{(i)} + b_k}} \right) - c_i \right] = \frac{\partial L}{\partial b_i}$$

Frobenius Norm

$$\|X\|_F^2 = \text{tr}(X^T X)$$

Training data: learn model parameters
Testing data: scoring the model

Validation:
See: select optimal hyperparameters
Test set: used as plug-in to see how well the architecture and hyperparameters would generalize to new data.

Test set used once after cross-validation (each fold is used K times for training/validation)
Early stopping prevents overfitting

Derivatives & Traces

$$\nabla_{\theta}(\Theta^T \chi) = \Theta$$

$$\nabla_x(x^T A x) = (A + A^T)x$$

$$\frac{d}{d\theta} \left(\frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T \hat{x}^{(i)})^2 \right), \theta = (X^T X)^{-1} (X^T Y)$$

$$\text{tr}(A) = \text{tr}(A^T)$$

$$\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$$

$$\nabla_X(Wx) = W^T$$

$$\frac{\partial \delta}{\partial W} = \frac{\partial \delta}{\partial X} X^T \quad Y = WX$$

$$\frac{\partial \delta}{\partial X} = W^T \frac{\partial \delta}{\partial Y}$$

$$-\text{tr}(W \sum_{i=1}^K x_i u_i^T) = -\text{tr}(WXY^T) \quad \checkmark HW \# 1$$

$$= \frac{\partial}{\partial W} \left(-\frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2 \right)$$

$$= -\frac{\partial}{\partial W} (\text{tr}(X^T W^T Y))$$

$$\nabla_{\theta} \Theta^T g = g$$

$$\nabla_{\theta} (\Theta^T H \Theta) = (H + H^T) \Theta$$

Linear Algebra

$$\textcircled{1} \quad Av_i = \lambda_i v_i \rightarrow \|v_i\|^2 = |\lambda_i|^2 \|v_i\|^2 \quad \text{norm of eigenvalues} = 1$$

$$\textcircled{2} \quad v_i^T A v_i = \lambda_i v_i^T v_i = \lambda_i \|v_i\|^2, \forall A \text{ b.d.} \Rightarrow \forall v_i, \lambda_i > 0$$

$$\textcircled{3} \quad A = U \Sigma V^T$$

↑ singular
eigenvalues of $A^T A$
 $A^T A$ is diagonal

$$AA^T = U \Sigma \Sigma^T V^T V = U \Sigma^2 V^T \quad \checkmark V \text{ orthogonal}$$

$$\textcircled{4} \quad Q^T Q = Q Q^T = I$$

$$\textcircled{5} \quad \|QV\| = \|V\|$$

$$\textcircled{6} \quad \det(Q^T Q) = \det(Q) = \det(I)$$

$$\text{cov}(Ax+b) = A \text{cov}(x) A^T$$

Arg, Det \equiv
represents the # of samples that belong to class j

$$\nabla_{w_i} L = \nabla_{w_i} \left(\frac{1}{m} \sum_{j=1}^m \left(\log \left(\frac{e^{w_j^T x^{(i)} + b_j}}{\sum_{k=1}^m e^{w_k^T x^{(i)} + b_k}} \right) - (w_j^T x^{(i)} + b_j) \right) \right)$$

$$= \frac{1}{m} \sum_{j=1}^m \left(\nabla_{w_i} \left(\log \left(\frac{e^{w_j^T x^{(i)} + b_j}}{\sum_{k=1}^m e^{w_k^T x^{(i)} + b_k}} \right) - (w_j^T x^{(i)} + b_j) \right) \right)$$

channel

$$\nabla_{w_i} \left(\log \left(\frac{e^{w_j^T x^{(i)} + b_j}}{\sum_{k=1}^m e^{w_k^T x^{(i)} + b_k}} \right) \right) = \frac{\sum_{k=1}^m \nabla_{w_k} (e^{w_k^T x^{(i)} + b_k})}{\sum_{k=1}^m e^{w_k^T x^{(i)} + b_k}}$$

$$\therefore \frac{1}{m} \left[\sum_{j=1}^m \left(\frac{e^{w_j^T x^{(i)} + b_j}}{\sum_{k=1}^m e^{w_k^T x^{(i)} + b_k}} \right) - z^{(i)} \right] = \frac{\partial L}{\partial w_i}$$

of times $y^{(i)} = i$

- Hinge Loss
- $$\nabla_w \mathcal{L} = \frac{1}{K} \sum_{i=1}^K \sum_j \epsilon_i y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \mathbb{1}_{\{\mathbf{w}^T \mathbf{x}^{(i)} + b \leq 1\}} \odot (-y^{(i)} \mathbf{x}^{(i)})$$
- Gradient Descent
- don't use a numerical gradient: too slow, expensive
 - batch: use all m training set examples
 - Minibatch: use K examples
 - Stochastic: use 1 example

Neural Networks

- Activation is usually RELU

$$-\text{Maxout} = \max(w_1^T x + b_1, w_2^T x + b_2)$$

(log loss)

Output Activation

$$\text{MSE} = \frac{1}{2} \sum (y^{(i)} - O(z^{(i)}))^2 \quad \begin{matrix} \text{cross} \\ \text{entropy} \\ \text{loss} \end{matrix}$$

$$\text{CE} = -\sum [y^{(i)} \log O(z^{(i)}) + (1-y^{(i)}) \log (1-O(z^{(i)}))] \rightarrow \text{softmax}$$

If z is large, the gradient ≈ 0 , therefore no learning.
Usually Softmax

Activations - introduce nonlinearity

- without an activator, everything is linear and simplifies to

$$\mathcal{Z} = \mathbf{W}^T \mathbf{x} + b$$

$$-\text{Sigmoid: } \frac{1}{1+e^x} = \sigma. \frac{d\sigma}{dx} = \sigma(1-\sigma)$$

- around $x=0$, linear, differentiable everywhere

- saturates at extremes, 0 gradient, no learning

- could zigzag

$$-\text{Hyperbolic: } \tanh(x) = 2\sigma - 1, \frac{d\tanh(x)}{dx} = 1 - \tanh^2(x)$$

- similar props to sigmoid passes through D

- ReLU: $\max(x, 0)$

- converges faster, linear if active, no saturation at large values

- non-neg, zig-zags, no learning if activation is 0

- Softplus: $\log(1+e^x)$

- empirically worse than ReLU

- Leaky ReLU: $x^+(x > 0) + 0.1x^-(x < 0)$

- learning if $x < 0$, if α is hyperparam, PReLU

- ELU: $\max(\alpha(e^x - 1), x)$

- computation of exponential is expensive

Initialization

- if too small, $\sigma \approx 0.01$, decay $\rightarrow 0$, fast learning
- if too large, $\sigma \approx 1$, explode, tanh gradients are small (causes fluctuation)
- Xavier initialization

$$\text{Var}(h_i) = \text{Var}(h_o) = \text{Var}(f_{in} L) = \text{Var}(L) \text{Var}(in)$$

Variance of units across all layers must be the same

$$\text{Var}(h_i) = \text{Var}\left(\sum_{j=1}^m w_{ij} h_{i-1,j}\right), \text{Var}(wh) = E^2(w) \text{Var}(h_{i-1}) + E^2(h) \text{Var}(w) + \text{Var}(w) \text{Var}(h_{i-1})$$

$$\text{if } E[w] = E[h] = 0,$$

$$\text{Var}(h_i) = \text{Var}(h_{i-1}) \cdot \sum_{j=1}^n \text{Var}(w_{ij})$$

$$\text{Var}(w) \times N_{in} = 1, \text{Var}(w) = \frac{1}{N_{in}} \quad \begin{matrix} \text{# of neurons} \\ \text{in a layer} \end{matrix}$$

Both

$$\mathcal{U}(0, \frac{2}{n_{in} + n_{out}})$$

2015:

$$\text{Uniform: } \left(-\frac{\sqrt{6}}{n_{in} + n_{out}}, \frac{\sqrt{6}}{n_{in} + n_{out}}\right)$$

if ReLU

$$H(y_p) = -\sum_i y_i \log(p_i) = L$$

2.1

$$\Omega(\Theta) = \frac{1}{2} \mathbf{W}^T \mathbf{W} = \frac{1}{2} \|W\|_F^2$$

- shrink weights before gradient update

$$-\tilde{J} = J + \frac{\alpha}{2} W^T W$$

$$\nabla_w \tilde{J} = \nabla_w J + \alpha W$$

$$W = (1-\alpha)W - \epsilon \nabla_w J$$

~ causes weights to shrink

- W could be near 0 $\|W\|_F^2$

Regularization

sacrifices higher train-loss for lower testing loss.

Stopping early

- return model with lowest validation error, train until the error is unchanged

Full norm penalty

$$2\Omega(\Theta) + R(\Theta)$$

Dropout approximates bagging

2.1

$$-\Omega(\Theta) = \sum \|W_j\|$$

- good for feature selection because gradient will make some weights go to exactly 0

Instead of sparse representations
 $\|w^{(i)}\|_1 = \sum |w^{(i)}_j|$

dataset augmentation

Multi-task learning

transfer learning

ensemble methods - avg multiple models

- bagging: K datasets/models small batch size

introduces noise to model

Batchnorm

makes each layer have unit statistics

$$\hat{x}_i = \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}, \quad \mu_i = \frac{1}{m} \sum_{j=1}^m x_i^{(j)}, \quad \sigma_i^2 = \frac{1}{m} \sum_{j=1}^m (x_i^{(j)} - \mu_i)^2$$

Allows higher learning rates to be used.

Optimization

- stochastic gradient descent / batch

$$g = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} J(\theta)$$

$$\theta = \theta - \epsilon g$$

using m samples from training set

Batch normalization computational graph

Momentum

$V = 0$

while $\text{last_diff} > \text{tolerance}$:

$$V = \alpha V - \epsilon \nabla_{\theta} J(\theta)$$

$$\theta += V$$

pushes for local, shallow, flat, minimum ↗

Nesterov momentum

$\alpha = 0.9$

temp: P
 $P = \alpha P - \epsilon g$
 $\theta += P + \alpha(P - \text{temp})$
 Also finds local minimum

Adagrad

learning rate decreased by grad norm

Let a denote running sum of squares of gradient

$$a = a + g \otimes g$$

"if we've already stepped far in one direction, we'll step less in that direction in the future!"

RMS Prop

- augments adagrad by making the gradient accumulator an exponentially weighted moving average

"Let's consider recent history of grad norm squared, and place a little less emphasis on our grad norm so we don't shut off one direction too much"

$$a = \beta a + (1-\beta) g \otimes g \quad \beta = 0.97$$

$$\theta = \theta - \frac{\epsilon}{\sqrt{a} + \nu} \otimes g$$

RMS Prop + momentum

"Let's consider recent history of both gradient and grad norm squared."

Let's try and step more in direction of recent gradients, & still place emphasis on recent norm squared, but not too much regarding our grad norm squared

$$a = \beta a + (1-\beta) g \otimes g \quad \theta = \theta + V \quad (\text{can also use nesterov})$$

$$V = \alpha V - \frac{\epsilon}{\sqrt{a} + \nu} \otimes g$$

Adams

- momentum-like + Adagrad/RMS prop-like, 2 moments + gradient
- $V = B_1 V + (1 - B_1) g$
- bias correction amplifies 2nd moment
- $B_1 = 0.9, B_2 = 0.999$
- $V = B_1 V + (1 - B_1) g$ \checkmark momentum-like
- $\alpha = B_2 \alpha + (1 - B_2) g \log \checkmark$ gradient normalization
- $\hat{V} = \frac{1}{1 - B_1^t} V$ \checkmark bias correction
- $\tilde{\alpha} = \frac{1}{1 - B_2^t} \alpha$ \checkmark all in one step R (α to your horizon)
- $\Theta = \Theta - \frac{\epsilon}{\hat{V}} \odot V$

Bias Correction

- fixes 1st & 2nd moments
- attenuates tails of gradient descent. makes sure it starts off well
- first steps more effective

CNN

$$(W - W_{ft} + h - h_{ft}) \\ \text{padding} + \frac{(W - W_{ft} + 2 \cdot \text{pad})}{\text{Stride}} + 1, \frac{h - h_{ft} + 2 \cdot \text{pad}}{\text{Stride}} + 1)$$

Miscellaneous

- exploding gradients
- gradient clipping, $g = \frac{g}{\|g\|_1} \cdot \epsilon$
- vanishing gradient \Rightarrow architecture, regularization

$$W = U A V^T, W^T = U A^T V^T$$

e_j in vector U_i , $h_i \in Z_i$

Ensemble Methods

1. Train multiple different models
2. Average results at test time

K independent models, avg model error will decrease by $\frac{1}{K}$

expected mean error square error, if dependent, $\frac{1}{K} \mathbb{E}[e_i^2] + \frac{K-1}{K} \mathbb{E}(e_i e_j) = \mathbb{E}(e_i^2) + r_i j$; perfect correlation

Dataset Augmentation

- flip, scale, rotate, crop images to generate more data
- inject noise into network
- label smoothing
- $L(c) = \sum_{c=1}^C y_c \log(p_c)$
- subsampling
- brightness adjustments
- gaussian blurring

Bagging

- bootstrap aggregation
- method for regularization
- construct K datasets by setting a data size N , and drawing w/ replacement from the original dataset to get N samples, we do this K times
- train different models using K datasets

Transfer Learning

- freeze most parameters of original network
- explain variations in data

multitask learning

- train to perform multiple tasks
- multiple tasks share common factors to

Dropout

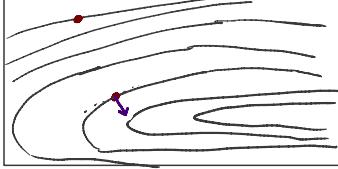
- approximates bagging procedure
- applies a mask of bernoulli r.v.'s
- each mask is like a different model
- N hidden units $\rightarrow 2^N$ model configurations
- each one good at predicting output
- dropout may cause units to encode different contexts
- dropout may cause units to encode redundant features

- minimizing $-\log L \rightarrow \max L$, min CE_{loss}
- ① Gradients computed using SGD will be noisier than gradient computed using Batch Gradient descent, because it is a single sample estimate
 - ② SGD will take more steps than Batch GD, but converges faster in terms of clock time
"Each step is super fast in SGD, and each step is slower in batch GD"
 - ③ SGD avoids the trap of poor local minima due to zig-zagging. Every sample is unlikely to have the same minimum as the full batch.
 - ④ Batch gradient descent is more expensive than SGD because it requires loading the entire data.

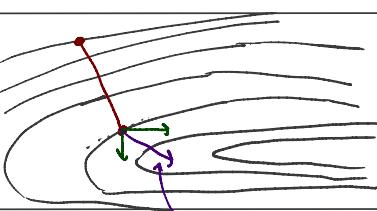
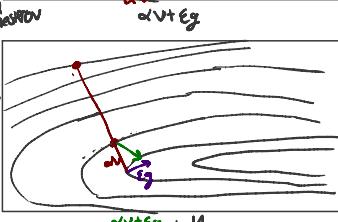
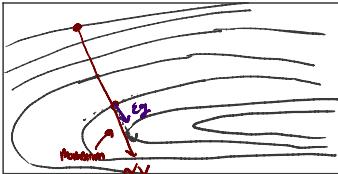
Difference between batchnorm in training VS. testing.

During training, it first computes the mean & variance of the mini-batch data in the unit-wise. BatchNorm normalizes the layer's input based on the mean & variance. It then applies γ, β for scaling and shifting (learnable parameters). Keeps track of the running averages for mean and variance through all batches. During testing, BatchNorm uses the fixed accumulated running averages from training for normalizing test data. The learned scaling/shifting γ, β are applied to normalized data.

Vanilla gradient descent is in the direction perpendicular to contour line



SGd of momentum



Adagrad

The update equations:

$$\alpha \leftarrow \alpha + g_t \\ \theta \leftarrow \theta + \frac{\epsilon}{\alpha} g_t$$

If gradient is accumulated more in a direction, then step $\frac{\epsilon}{\alpha}$ is less in that direction

If gradient accumulated is less in a direction. Then the step $\frac{\epsilon}{\alpha}$ is more in that direction. From plot, we can see that when it traversed from A to B, the gradients accumulated along W_2 direction is more and gradient along W_1 is small. Here, the step along W_2 is more and step across W_1 is less.

Adagrad converges faster to minima w/o much zig-zagging

Gradient descent + momentum scheme, find a general expression of V_t in terms of gradients g_1, g_2, \dots, g_t and ϵ , considering an initial value of momentum $V_0=0$.

$$V_t = \alpha V_{t-1} - \epsilon g_t \quad V_i = \alpha V_0 - \epsilon g_i = -\epsilon g_i$$

$$\theta_t = \theta_{t-1} + V_t$$

$$V_2 = \alpha V_1 - \epsilon g_2 = \alpha(-\epsilon g_1) - \epsilon g_2 = -\epsilon(g_1 + \alpha g_2)$$

$$V_3 = \alpha V_2 - \epsilon g_3 = \alpha(-\epsilon(g_1 + \alpha g_2)) - \epsilon g_3 = -\epsilon(g_1 + \alpha g_2 + \alpha^2 g_3)$$

$$= -\epsilon(g_1 + \alpha g_2 + \alpha^2 g_3 + \alpha^3 g_4)$$

$$\Rightarrow V_t = -\epsilon [g_t + \alpha g_{t-1} + \alpha^2 g_{t-2} + \dots + \alpha^{t-1} g_1]$$

$$\mathbb{E}[V_t] = \mathbb{E}[-\epsilon(g_t + \alpha g_{t-1} + \alpha^2 g_{t-2} + \dots + \alpha^{t-1} g_1)]$$

$$= -\epsilon [\mathbb{E}[g_t] + \alpha \mathbb{E}[g_{t-1}] + \dots + \alpha^{t-1} \mathbb{E}[g_1]]$$

$$\text{Consider that the gradients above are i.i.d. random variables with } \mathcal{N}(0, \sigma^2). \quad \mathbb{E}[g_t] = \frac{-\epsilon M}{1-\alpha} \sum_{j=1}^t (1-\alpha^j)$$

Find the expected value of weights θ_t at $t=3$.

$$\theta_1 = \theta_0 + V_1$$

$$\theta_3 = \theta_0 + V_1 + V_2 + V_3$$

$$= \theta_0 - \epsilon \mathbb{E}[g_1] - \epsilon (\mathbb{E}[g_2] + \alpha \mathbb{E}[g_1]) - \epsilon (\mathbb{E}[g_3] + \alpha \mathbb{E}[g_2] + \alpha^2 \mathbb{E}[g_1])$$

$$= \theta_0 - \epsilon M - \epsilon(\mu + \alpha \mu) - \epsilon(\mu + \alpha M + \alpha^2 \mu)$$

$$= \theta_0 - \epsilon M - \epsilon M - \epsilon \alpha M - \epsilon \mu - \epsilon \alpha M - \epsilon \alpha^2 M$$

$$= \theta_0 - 3\epsilon M - 2\epsilon \alpha M - \epsilon \alpha^2 M$$

$$= \theta_0 - \epsilon M(3 + 2\alpha + \alpha^2)$$

$$= \theta_0 - \epsilon M[1 + (1+\alpha) + (1+\alpha+\alpha^2)]$$

$$\mathbb{E}[V_t] = -\epsilon M \left[\frac{1-\alpha^t}{1-\alpha} \right]$$

$$\mathbb{E}[\theta_3] = \mathbb{E}[\theta_0] + \sum_{t=1}^3 \left[-\epsilon M \frac{1-\alpha^t}{1-\alpha} \right]$$

$$= \mathbb{E}[\theta_0] - \frac{\epsilon M}{1-\alpha} \left[\sum_{t=1}^3 (1-\alpha^t) \right]$$

$$= \mathbb{E}[\theta_0] - \frac{\epsilon M}{1-\alpha} \left[3 - \left(\frac{1-\alpha^3}{1-\alpha} \right) \right]$$

$$\|x\|_\infty = \max_i |x_i|$$

$$\text{LSE}(x) = \ln \left(\sum_{i=1}^n e^{x_i} \right)$$

$$\|x\|_\infty \leq \text{LSE}(x) \leq \|x\|_\infty + \ln(n)$$

$$P = \max_i \{ |x_1|, \dots, |x_n| \}$$

$$e^P \leq e^{|x_1|} + e^{|x_2|} + \dots + e^{|x_n|}$$

$$e^P \leq \sum_{i=1}^n e^{|x_i|}$$

$$\sum_{i=1}^n e^{|x_i|} \leq n e^P \quad \begin{matrix} \text{mean don't grow} \\ \text{every time} \end{matrix}$$

$$\ln(e^P) \leq \ln \left(\sum_{i=1}^n e^{|x_i|} \right) \leq \ln(n e^P)$$

$$P \leq \ln \left(\sum_{i=1}^n e^{|x_i|} \right) \leq \ln(n) + P$$

Under what condition on x will

$$\theta_t = \theta_0 - \epsilon \sum_{j=1}^t \sum_{i=1}^{j-1} \alpha^{j-i} g_i$$

$$\mathbb{E}[\theta_t] = -\epsilon \sum_{j=1}^t \sum_{i=1}^{j-1} \alpha^{j-i} \mathbb{E}[g_i]$$

$$= -\epsilon M \sum_{j=1}^t \frac{1-\alpha^j}{1-\alpha}$$

$$\mathbb{E}[\theta_t] = \frac{\epsilon M}{1-\alpha} \left(t - \frac{\alpha(1-\alpha^t)}{1-\alpha} \right)$$

$$\|x\|_\infty \leq \ln \left(\sum_{i=1}^n e^{|x_i|} \right) \leq \|x\|_\infty + \ln(n)$$

$$\|tx\|_\infty \leq \ln \left(\sum_{i=1}^n e^{|tx_i|} \right) \leq \|tx\|_\infty + \ln(n)$$

$$t\|x\|_\infty \leq \ln \left(\sum_{i=1}^n e^{|tx_i|} \right) \leq \|x\|_\infty + \ln(n)$$

$$\|x\|_\infty \leq \frac{1}{t} \ln \left(\sum_{i=1}^n e^{|x_i|} \right) \leq \|x\|_\infty + \frac{\ln(n)}{t}$$

$$\|x\|_\infty \leq \text{LSE}(x) \leq \|x\|_\infty + \ln(n)$$

$$\ln \left(\sum_{i=1}^n e^{|x_i|} \right) \geq \|x\|_\infty$$

$$\sum_{i=1}^n e^{|x_i|} > e^P$$

$$\ln \left(\sum_{i=1}^n e^{|x_i|} \right) > P$$

$$\text{LSE}(x) > \|x\|_\infty$$

$$\ln \left(\sum_{i=1}^n e^{|x_i|} \right) = P + \ln(n)$$

$$\sum_{i=1}^n e^{|x_i|} = n e^P \quad \begin{matrix} \text{every } |x_i| = \\ \text{for } j \neq i \end{matrix}$$

K-nearest neighbors

$$\text{error} = \frac{\#\text{num wrong}}{\#\text{total}}$$

In order to minimize classification error on this test set, we don't necessarily need to choose a value of K that minimizes training, even if sampled from same distribution.

Small value of K \rightarrow overfitting

Large value of K \rightarrow underfitting | reduces effect of outliers

Use cross validation to choose K.

Activation functions where vanishing gradient is a problem.

- if it has zero gradient somewhere, then there's a problem

- ReLU, O, tanh

MORE ON BATCH NORM

- adds noise/regularization because

convoluted our minibatches, not entire dataset

each delta distribution has a little bit of noise

- batch normalization accelerates training by requiring fewer iterations to converge to a given loss value

- batch norm has learnable parameters

- linear transformation @ test time

- introduces noise to hidden layer's activation because the mean and std are estimated w/ a mini-batch

- use running averages at test time

Adam combines effects of Adagrad & momentum

Learning rate decay may cause sudden drops in loss landscape

WEIGTS LEARNED IN L1 REGULARIZATION WILL BECOME more sparse than the weights learned using L_2 regularization

ENSEMBLE methods: if NN's make independent errors, then they are likely to make errors on diff trials.

We only get a trial wrong when several neural networks make the same errors. Here, lower correlation \mapsto better accuracy.

$$v_t = -\epsilon \sum_{i=1}^t \alpha^{t-i} g_i$$

$$\begin{aligned} \text{Cov}(v_t) &= E[v_t v_t^\top] \\ &= \epsilon^2 E\left[\sum_{i=1}^t \alpha^{t-i} g_i \left[\sum_{j=1}^t \alpha^{t-j} g_j\right]^\top\right] \\ &= \epsilon^2 E\left[\sum_{i=1}^t \sum_{j=1}^t \alpha^{t-i} \alpha^{t-j} g_i g_j^\top\right] \\ &= \epsilon^2 E\left[\sum_{i=1}^t \alpha^{2(t-i)} \sigma^2\right] = \epsilon^2 \sigma^2 \left(\frac{1-\alpha^2}{1-\alpha}\right) I \end{aligned}$$

$$\begin{aligned} E[V_t] &= (1-\beta_1) \sum_{i=1}^t \beta_1^{t-i} E[g_i] \\ &= (1-\beta_1) M \sum_{i=1}^t \beta_1^{t-i} = M(1-\beta_1) \left[\frac{(1-\beta_1^t)}{1-(1-\beta_1)} \right] \end{aligned}$$

$$E[V_t] = M(1-\beta_1^t)$$

$$\text{Var}[V_t] = M \frac{(1-\beta_1^t)}{1-\beta_1^t} = M$$

$$\sum_{k=m}^n ar^k = \frac{a(r^m - r^{n+1})}{1-r}$$

(d) (5 points) Consider a convolutional neural network (CNN) with N conv-pool layers.

Each conv-pool layer uses five filters of size 3×3 , stride = 1 and padding = 1, followed by $\text{relu}()$, followed by 2×2 max pooling with stride 2. The input image is of size $256 \times 256 \times 3$. What is the number of trainable parameters (including biases) in the first conv-pool layer?

$$\begin{aligned} & \text{- trainable parameters } (\text{filter width} \times \text{filter height} \times 3) + 1) \times (\# \text{of filters}) \\ & = 28 \times 5 = 140 \end{aligned}$$

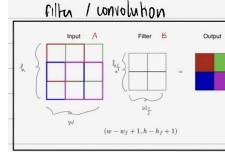
ratio of class samples is 10, pick $\alpha = 1, \beta = 10$

$$\frac{\text{Weight for class}_i}{\# \text{num_samples}} = \frac{\text{total samples}}{\# \text{num_classes}}$$

$$\frac{2}{n} \sum_{i=1}^n \frac{\partial}{\partial W} (\max_{1 \leq k \leq p} \sum_{j=1}^m |W_{kij}|)$$

$$\frac{2}{n} \sum_i (|W_{j,1}| + |W_{j,2}| + \dots + |W_{j,p}|)$$

$$\frac{\partial q}{\partial W^{(i)}} = \left[\frac{\partial q}{\partial W^{(i)}_{11}}, \frac{\partial q}{\partial W^{(i)}_{12}}, \dots, \frac{\partial q}{\partial W^{(i)}_{1p}} \right]$$



a VALID convolution: A and B fully overlapping

each layer has multiple filters, let us call that number n_f , each of which has its own slice outputs put together into tensor: $(w-w+1, h-h+1, n_f) \rightarrow$ passed into PReLU $\rightarrow h_i$, passed on

sparsely connected layers

- 1 output neuron only connected to values inside the convolution
- e.g. above, each neuron only connected to $5 \times 5 \times 3$

spars interaction means there should be more layers

e.g. 6 filters are $(4 \times 4 \times 3)$, so params: $(4 \times 4 \times 3 + 1) \times 6 = 196$ params
CNN has $(227 \times 32 \times 3 + 1) \times 6$ neurons in layer \approx millions

use padding: $(w-w+1 + 2\text{pad}, h-h+1 + 2\text{pad})$

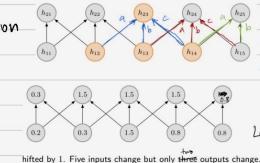
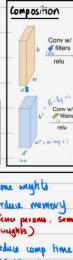
use stride: $(\frac{w-w+1 + 2\text{pad}}{\text{stride}}, \frac{h-h+1 + 2\text{pad}}{\text{stride}})$

use pooling: $(\frac{w-w+1}{\text{stride}}, \frac{h-h+1}{\text{stride}})$ \rightarrow no additional params
small spatial invariance

b) sizing formulas for output size

input: $h=32, w=32, d=3$ } same depth, 3D $\rightarrow (w-w+1, h-h+1)$
filter: $w=5, h=5, d=3, s=2$ } stride = 2, $d=3$ $\rightarrow (w-w+1, h-h+1, n_f)$
note: $(5 \times 5 \times 3) + b = 75 + 1 = 76$ params

each output matrix is a slice



Some weights
(few neurons, some weights)
reduce comp time
 $O(nk)$ instead of $O(n^2)$

C1 contains six 5×5 conv filters. stride = 1, pad = 0
Size of feature maps at C1? $32 - 5 + 1 = 28 \rightarrow (28 \times 28 \times 6)$
Number of parameters in C1 layer? $(5 \times 5 + 1) \times 6 = 156$ params

S2 is a 2×2 pooling layer applied at stride 2.

Size of feature maps at S2? $(14 \times 14 \times 6)$

Number of parameters in S2 layer? 0

C3 contains sixteen 3×3 conv filters.

Size of feature maps at C3? $(7 \times 7 \times 16)$

Number of parameters in C3? 0

- sizing: width, height from above, H of filters

- trainable params: $(\text{filter width} \times \text{filter height} \times 3) + 1) \times \# \text{of filters}$

- pooling has 0 parameters

(5) VGGnet

- CONV uniform: 3×3 , stride 1, pad 1 } 1171 \rightarrow 731,

- POOL uniform: 2×2 max, stride 2

- smaller fields need more layers, less params, but more nonlinearity
non-linearity \leftarrow but $(\text{width} \times \text{height}) \times 3 \times (\text{width} \times \text{height})$, more than $3 \times$ activations

- 138M params, concentrated in FC layers (122M)

(6) GoogleNet = deeper Inception module, no FC, 5M params

- previous layers \rightarrow $\frac{1}{2} \times \frac{1}{2}$ \rightarrow concatenate \rightarrow next layer
max pool 2×2

① LeNet 5

- neuron $(5 \times 5 \times 1)$ } 122,304
- # of neurons $(28 \times 28 \times 6)$

1. [Model] CONV 6 convolutional filters, filter size applied at stride 1.
2. [ReLU] ReLU activation function, takes the max of zero and the input, then passes through sigmoid.
3. [MaxPool] Max pooling layer, takes the max of the four values at each position.
4. [Conv] CONV 16 convolutional filters, filter size applied at stride 1.
5. [ReLU] ReLU activation function, takes the max of zero and the input, then passes through sigmoid.
6. [FC] FC 1000 fully connected layer with 1000 units.
7. [Out] OUT: Softmax activation for each digit.

② AlexNet

- input processing
- use pooling with overlapping
- augmentation
- extract patches + inflation, 1.5% error
- scale PCs of color = 1%, error
- dropout - p = 0.5

1. [Input] Input image, $227 \times 227 \times 3$.
2. [Dropout] POOL 3x3 win, stride 2.
3. [Conv] CONV 11x11, stride 4, padding 5.
4. [Norm] NORM 5x5 kernel, stride 1, padding 2.
5. [ReLU] ReLU activation function, takes the max of zero and the input, then passes through sigmoid.
6. [Dropout] POOL 3x3 win, stride 2.
7. [Conv] CONV 5x5 kernel, stride 1, padding 2.
8. [Norm] NORM 5x5 kernel, stride 1, padding 2.
9. [ReLU] ReLU activation function, takes the max of zero and the input, then passes through sigmoid.
10. [Dropout] POOL 3x3 win, stride 2.
11. [Conv] CONV 3x3 kernel, stride 1, padding 1.
12. [Norm] NORM 3x3 kernel, stride 1, padding 1.
13. [ReLU] ReLU activation function, takes the max of zero and the input, then passes through sigmoid.
14. [FC] FC 1000 fully connected layer with 1000 units (cross entropy loss).
15. [Out] OUT: Softmax layer.

③ zinet with different hyperparameters

- e.g. 7x7 had more freq coverage than 11x11
- they increased filters: 384 to 512/1024
dropped from 164 to 117