

Key problem in ML: How to generalize on unseen data?

Demo.

One technique: Regularization — any modification to learning algorithm that is intended to reduce its generalization (test) error but not its training error.

There are many ways ; Restrict the model- hyper parameters: no of nodes in a decision tree.

add extra terms in the objective fn itself.

Combine multiple hypotheses - Ensemble methods.

III. Restrict the parameter space. — most common.

$$\text{SVM} = \gamma w^T w + C \sum |e_i|$$

We want — a good trade where bias is increased for reduced variance.

Bayesian viewpoint — Adding prior to the estimation : eg: L_2 reg. for linear regression is \sim MAP estimate with Normal prior on w .

In the context of DL, we regularize to generalize over true data-generating distributions that are often very complicated

Eg^o image classification task - true data generating distro is the entire universe (?)
 /speech

Only weights are regularized - not the biases.

reasons - 1) Biases only control single variable (through b^d)
 (weights quantify the interactions b/w features)

2) Regularizing bias induces overfitting
 ∵ intuitively b^d for neuron firing shouldn't be restricted.

Parameter norm penalties:

$$\hat{J}(\theta; x, y) = J(\theta; x, y) + \alpha \Omega(\theta).$$

Bias-variance Dilemma:

Error of a function approximator,

$$E = \lim_{N \rightarrow \infty} \frac{1}{2N} \sum_{i=1}^N \sum_k \left\{ f_k(x_i; w) - y_k^i \right\}$$

$f_k \rightarrow$ estimate for k^{th} output variable
Then $y \in \mathbb{R}^k$

$$= \frac{1}{2} \sum_k \int \left\{ f_k(x; w) - y_k^* \right\} P(y_k|x) dy_k^* dx$$

Define $\langle y_k^* | x \rangle = \int y_k^* P(y_k|x) dy_k^*$

$$\langle (y_k^*)^2 | x \rangle = \int (y_k^*)^2 P(y_k|x) dy_k^*$$

Can show $E = \frac{1}{2} \sum_k \left\{ f_k(x; w) - \langle y_k^* | x \rangle \right\}^2 P(x) dx$

$$+ \frac{1}{2} \sum_k \left\{ \langle y_k^* | x \rangle - \langle y_k | x \rangle^2 \right\} P(x) dx$$

Minimal Error occurs when.

$$f_k(x_k; w) = \langle y_k^* | x \rangle$$

Can't change this.

However, since we don't know $\langle y_k | x \rangle$
we have to tweak w_3 in f to make error
as small as possible. ($\hat{y}(x) \equiv f(x, w)$)

$$\text{Simply put, } E = \frac{1}{2} \int \{ \hat{y}(x) - \langle y | x \rangle \}^2 p(x) dx \\ + \frac{1}{2} \int \{ \langle y^2 | x \rangle - \langle y | x \rangle \}^2 p(x) dx.$$

A measure of how close the approximator to the
desired one is given by

$$\{ \hat{y}(x) - \langle y | x \rangle \}^2 \quad (1)$$

Particular

This will depend on the dataset where the
network was trained.

To eliminate this dependence, Let's average (1)
over all possible datasets.

$$\text{measure of error} = \mathbb{E}_D [\{ \hat{y}(x) - \langle y | x \rangle \}^2] \quad (\text{f' of } x).$$

$\mathbb{E}_D \rightarrow$ Expectation over all datasets.

(2)

A perfect predictor will make it zero.

$$[\hat{y}(x) - \langle y|x \rangle] = \left\{ \hat{y}(x) - \mathbb{E}_D[\hat{y}(x)] \right\} + \mathbb{E}_D[\hat{y}(x) - \langle y|x \rangle]$$

$$\hat{E} = \mathbb{E}_D[\{\hat{y}(x) - \langle y|x \rangle\}^2] = (\text{can be shown to be equal to below})$$

$$\left\{ \mathbb{E}_D[\hat{y}(x) - \langle y|x \rangle]^2 \right\} + \mathbb{E}_D[\{\hat{y}(x) - \mathbb{E}_D[\hat{y}(x)]\}^2]$$

$$\text{Bias}^2 \quad \text{Variance}$$

$$\left\{ \mathbb{E}_D[\hat{y}(x) - \langle y|x \rangle]^2 \right\} + \mathbb{E}_D[\{\hat{y}(x) - \mathbb{E}_D[\hat{y}(x)]\}^2]$$

$$\text{Total Error} = \text{Bias}^2 + \text{Variance}$$

Bias — The extent to which the average (over all of the network f^n differs from the desired function.)

Variance — measures the extent to which the network $f^n \hat{y}(x)$ is sensitive to the choice of a particular dataset D.

(3)

Have to reduce both bias & variance

62

Now, Bias & variance are functions of x .

Lets look at Bias_{avg} , $\text{Variance}_{\text{avg}}$ by integrating out the x .

$$\therefore \text{Bias}_{\text{avg}}^2 = \frac{1}{2} \int \left\{ E_D [\bar{y}(x) - \langle y|x \rangle]^2 P(x) dx \right.$$

$$\text{Variance}_{\text{avg}} = \frac{1}{2} \int E_D \left[\left\{ \bar{y}(x) - E_D [\bar{y}(x)] \right\}^2 P(x) dx \right]$$

Lets study the behaviour of B & V .

Consider a regression problem with the target function being a smooth parabola.

Data is generated by adding a small noise to $h(x)$.

$$\hat{y} = h(x) + \epsilon$$

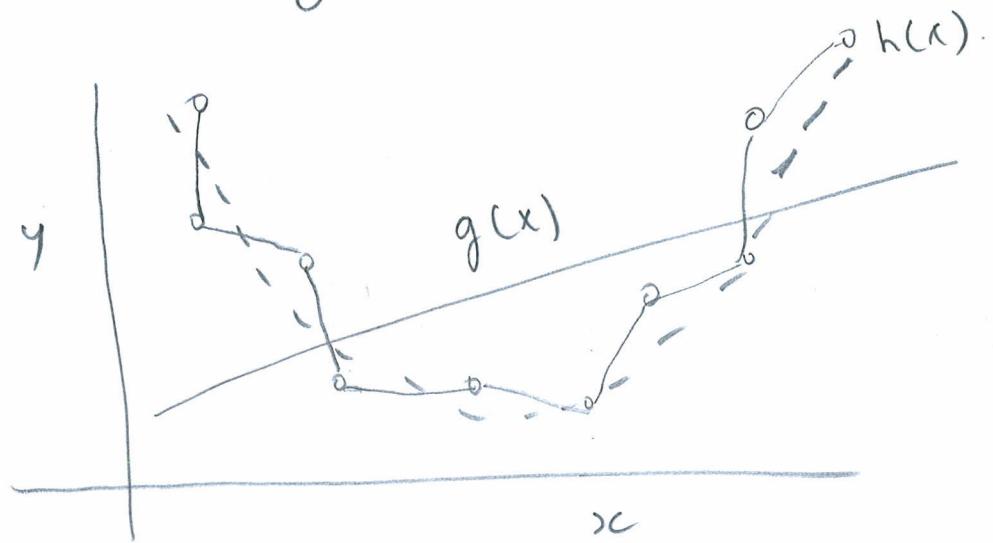
Optimal mapping function $\langle y|x \rangle = h(x)$.

Now, Lets choose $\bar{y}(x)$ to be some $g(x)$ which is completely ill of data.

In this case the variance term in E_m vanishes because $E_D [\bar{y}(x)] = g(x) = \bar{y}(x)$ [Ild of data].

(4)

However the bias term will be typically high unless we have some prior knowledge that helps us to choose $g(x)$.



The opposite extreme is to choose

$$\hat{y}(x) = h(x^n) + \epsilon \text{ & interpolate}$$

$$\begin{aligned} \text{In this case, Bias} &= 0 \quad \therefore E_D[\hat{y}(x)] = E_D[h(x) + \epsilon] \\ &= h(x) = \langle y | x \rangle. \end{aligned}$$

However, Variance

$$= E_D[(\hat{y}(x) - E_D(\hat{y}(x)))^2]$$

$$= E_D[(\hat{y}(x) - h(x))^2] = E_D[\epsilon^2]$$

= Variance of noise. (high).

Therefore, there is a trade-off b/w BEV.

f' that closely fits the data - how bias / high variance
 [overfitting] \Rightarrow large Em.

Variance can be decreased by smoothing the f'

But taken too far - bias becomes too large
 \Rightarrow large Em. [underfitting]

Question: How to min both BEV?

One answer: Have a large (infinite) dataset.

\Rightarrow use complex models & reduce bias.
 see enough samples to reduce variance.

One reason why DNN works well on large-scale datasets.

Question: How do you do that on limited data.
 [which is often the case in practice].

Another way: use some prior knowledge that you have about unknown $h(x) \rightarrow$ ^{example of} ^{sq waves.}

Now, use that knowledge to constrain $\hat{y}(x)$ so that bias is not increased too much.
(Biased models make better decisions).

This is known as regularization: Any modification to learning algo so that generalization is improved.

Regularization: Several strategies. (Go to page 51)

Strategy 1: Less is more! [Parameter norm penalties].

Data that we get with AI tasks are often sparse. [Demo of images from white noise].

When data is sparse, the dominant component of error.

Cognitive systems performs well with sparse data.

Combating bias \rightarrow require rich class of models
~~complement~~ $\nu \rightarrow$ placing restrictions on these models (7).

Examples of simple models (heuristics) performing remarkably better (Harry Markowitz - mean-variance portfolio allocation) vs $\frac{1}{N}$ portfolio.

\Rightarrow Simpler models (high bias) are better for sparse data.

Summary: Look for simpler models to improve \hat{E}_m .

Parameter norm penalties:

$$\hat{\mathcal{J}}(\theta; x, y) = \mathcal{J}(\theta; x, y) + \alpha \cdot R(\theta).$$

[Add a penalty term constraining the parameter space to the objective f^n .]

[DNN-biases are not regularized, go to pg. 57]

29/1

L^2 regularization (weight decay, ridge, Tikhonov).

$\mathcal{J}(\theta) = \frac{1}{2} \|\omega\|_F^2 \rightarrow$ Drives weight closer
decay factor to origin (can be anything)

$$\hat{\mathcal{J}} = \frac{\alpha}{2} \omega^\top \omega + J \quad \text{but Origin is chosen : tve & ve unknown}$$

$$\nabla_{\omega} \hat{\mathcal{J}} = \alpha \omega + \nabla_{\omega} J \quad \text{Forces some neurons to behave linearly.}$$

$$\omega^{k+1} = \omega^k - \epsilon (\alpha \omega^k + \nabla_{\omega} J)$$

$$= (1 - \epsilon \alpha) \omega - \epsilon \nabla_{\omega} J$$

Shrink the weight by a const factor on each step before gradient update

Let's study what L^2 does in the entire course of training.

P.T.O.

(g)

Let $\omega^* = \arg \min_{\omega} J(\omega) \rightarrow$ minimizer for un-regularized cost.

As before, let's make a quadratic approximation of J around ω^* .

$$\hat{J}(\theta) = J(\omega^*) + \frac{1}{2} (\omega - \omega^*)^T H (\omega - \omega^*)$$

Gradient term is absent $\because \omega^*$ is a minimizer.

$$\nabla_{\omega} \hat{J}(\omega) = H(\omega - \omega^*)$$

minima occurs when $\nabla_{\omega} \hat{J}(\omega) = 0$.

with regularization,

$$\alpha \hat{\omega} + H(\hat{\omega} - \omega^*) = 0$$

$\hat{\omega}$ - minima with regularized weight.

$$(H + \alpha I) \hat{\omega} = H \omega^*$$

$$\hat{\omega} = (H + \alpha I)^{-1} H \omega^*$$

as $\alpha \rightarrow 0$, $\hat{\omega} \rightarrow \omega^*$ (obvious)

But what happens with growing α ?

As usual, since H is real & symmetric

$$\exists \Phi, \Lambda, \text{ s.t. } H = \Phi \Lambda \Phi^T \text{ (E.V.D.)}$$

$$\therefore \hat{w} = (\Phi^\top \Phi + \alpha I)^{-1} \Phi^\top w^*$$

$$= [\Phi (\Lambda + \alpha I) \Phi^\top]^{-1} \Phi^\top \Lambda \Phi^\top w^*$$

$$= \Phi (\Lambda + \alpha I)^{-1} \Lambda \Phi^\top w^*$$

\Rightarrow weight decay rescales w^* along the axes defined by Evec of H .

Specifically, \hat{w} aligned with i^{th} Evec, is rescaled by a factor of $\frac{\lambda_i}{\lambda_i + \alpha}$.

\Rightarrow Along the directions where H has relatively higher Evals, $\lambda_i \gg \alpha$, the effect of L^2 is small. when $\lambda_i \ll \alpha$, weight gets shrunk (Demo).

Example with linear regression.

$$J = \|w^\top x - y\|_2^2 + \frac{1}{2} \alpha \|w\|_2^2$$

$$\Rightarrow w^* = (x^\top x)^{-1} x^\top y$$

$$\hat{w} = (x^\top x + \alpha I)^{-1} x^\top y$$

(11)

Now, $X^T X \rightarrow$ Covariance matrix.

* L^2 increases the "perceived" variance.

\Rightarrow weights on features whose covariance with the target output is low w.r.t to the new variance) are shrunked.

L^2 is also equivalent to MAP estimation.
with gaussian priors

L^2 is not sparse - it still makes weights to have no-zero values.

Desirable to have zero-values for weights that are unused.

solution: L^1 regularization.

$$\Omega(\theta) = \|\omega\|_1 = \sum_i |\omega_i|$$

$$\tilde{J}(\omega; x, y) = \alpha \|\omega\|_1 + J(\omega; x, y).$$

29/1

71

$$\nabla_{\omega} \tilde{J}(\omega; x, y) = \alpha \text{sign}(\omega) + \nabla_{\omega} J(x, y; \omega).$$

Again making quadratic approx,

$$\nabla_{\omega} \tilde{J}(\omega) = H(\omega - \omega^*).$$

Since L' does not have clean algebraic expression in case of full H , lets further simplify by assuming H is diagonal.

$$\text{i.e., } H = \text{diag}([H_{1,1}, \dots, H_{n,n}]), \quad H_{ii} > 0$$

[Easy to get with PCAed data].

$$\text{Now, } \tilde{J}(\omega; x, y) = J(\omega^*; x, y) +$$

$$\sum_i \left[\frac{1}{2} H_{ii} (\omega_i - \omega_i^*)^2 + \alpha |\omega_i| \right]$$

-(1)

Let $\tilde{\omega}_i$ be the minimizer with L' .

Have to find minimizer of (1) w.r.t. ω_i

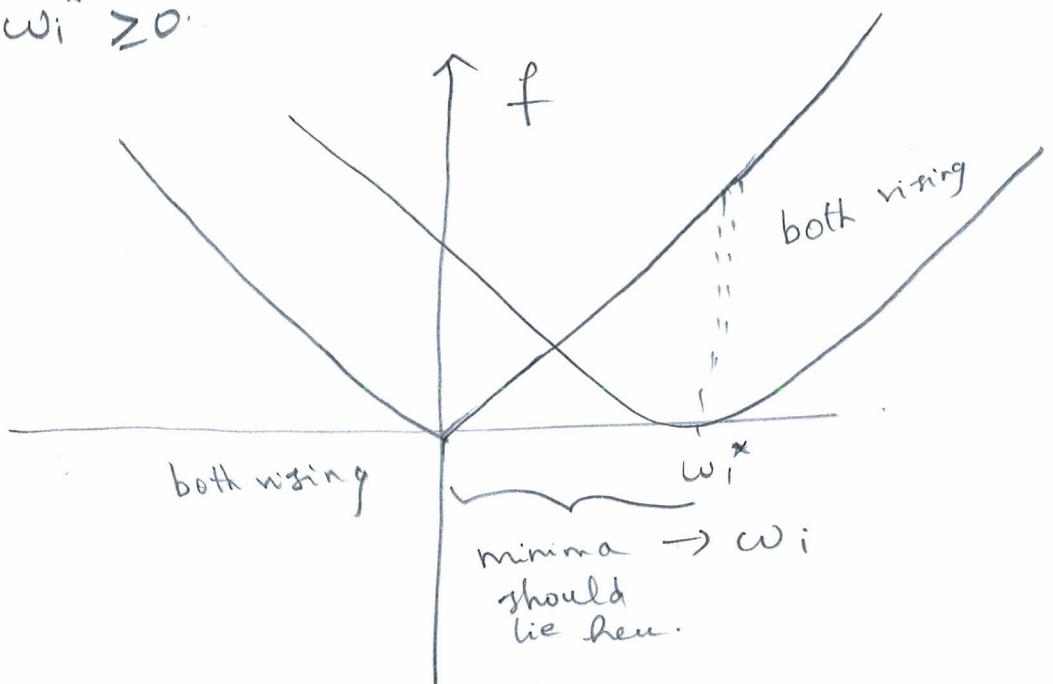
First term does not contribute to minima 'cause it is α of ω_i . Lets concentrate on the 2nd term.

29/11 We have to find the minima of the following. 72

$$\frac{1}{2} H_{ii} (w_i - w_i^*)^2 + \alpha |w_i| = f$$

↓ ↑
Parabola V-shaped curve.

Assume $w_i^* \geq 0$.



Also, algebraically, if $w_i^* \geq 0$,

$$\underbrace{w_i}_{\text{if } w_i \text{ is positive}} \quad \text{and} \quad \underbrace{\text{If }}_{\text{if } w_i \text{ is negative}} \quad \nabla_{w_i} f = (w_i - w_i^*) H_{ii} + \alpha \text{sign}(w_i)$$

for critical point, $\nabla_{w_i} f = 0$.

if $w_i^* > 0$, w_i can't be negative \therefore

$\nabla_{w_i} f$ should be zero.

$$\text{Further } \text{sign}(w_i) = w_i \Rightarrow \tilde{w}_i = w_i^* - \frac{\alpha}{H_{ii}}$$

Also from above graph, $0 \leq \tilde{w}_i \leq w_i^*$

$$\Rightarrow \tilde{w}_i^* = \max \left\{ w_i^* - \frac{\alpha}{H_{ii}}, 0 \right\}$$

29/1

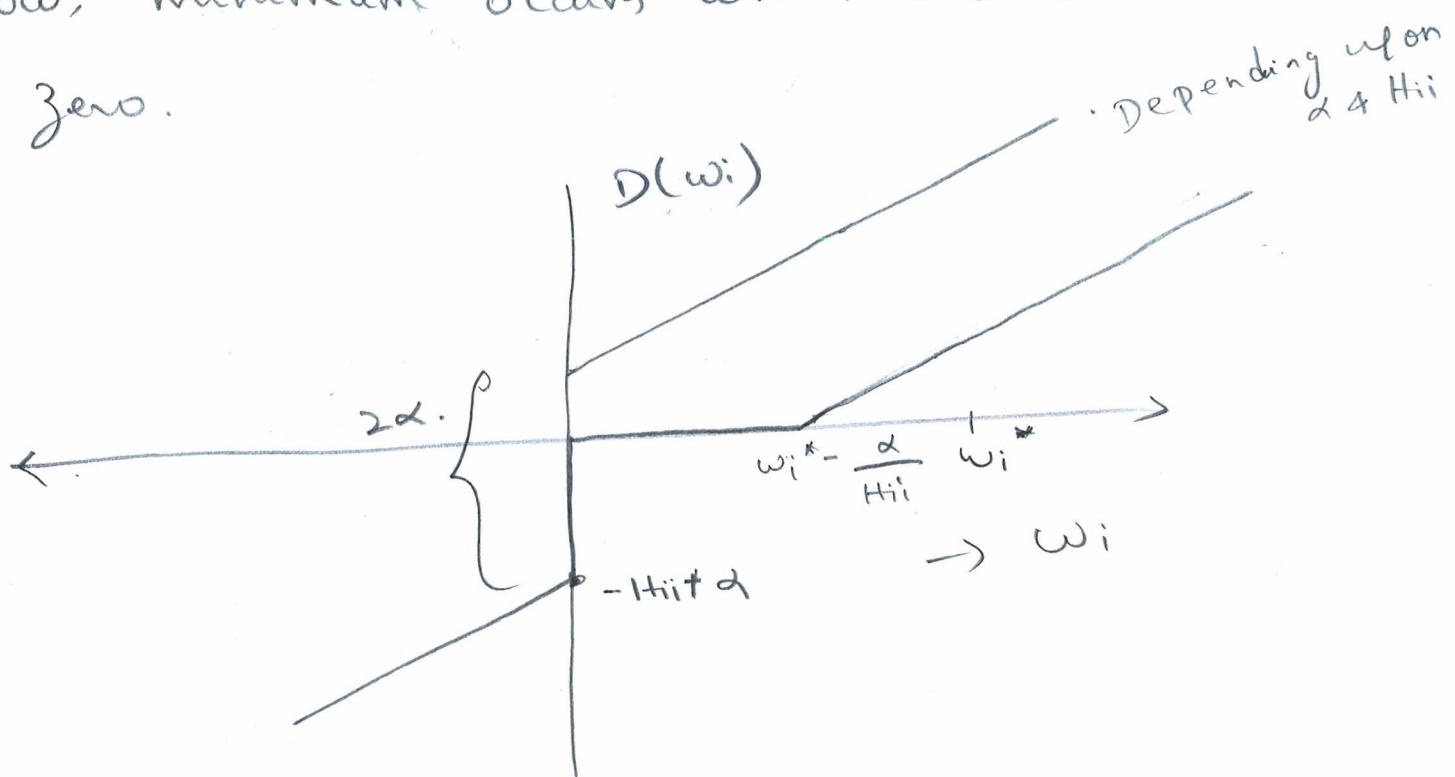
13

Another argument with derivative. Let $w_i \geq 0$.

$$D(w_i) = H_{ii} (w_i - w_i^*) + \alpha \text{sign}(w_i)$$

$\frac{\partial D(w_i)}{\partial w_i}$

Now, minimum occurs when $D(w_i)$ crosses zero.



\therefore zero crossing occurs either at 0, or at $w_i^* - \frac{\alpha}{H_{ii}}$

$$\text{Thus, } \hat{w}_i^* = \max \left\{ w_i^* - \frac{\alpha}{H_{ii}}, 0 \right\}.$$

A similar argument can be made when $w_i < 0$.

$$\text{Thus, } \hat{w}_i = \text{sign}(w_i^*) \max \left\{ |w_i^*| - \frac{\alpha}{H_{ii}}, 0 \right\}$$

(15).

74 now, let's interpret this.

For the case when $w_i^* > 0$, there are two cases.

i) $\alpha w_i^* \leq \frac{\alpha}{H_{ii}} \Rightarrow \tilde{w}_i^* = 0$

\Rightarrow along direction i , regularization overwhelms the other.

H_{ii} is small $\Rightarrow f^n$ varies less along that direction specified by i

(contribution of J to \tilde{J} is overwhelmed by L').

ii) when $w_i^* > \frac{\alpha}{H_{ii}}$,

L' just shifts the minimum in the direction i by $\frac{\alpha}{H_{ii}}$.

Summary : L' induces sparsity & might be better in some cases where there are too-many parameters [Again simplifying the model to increase the Bias]. Demo.

Info: L' is MAP with isotropic laplacian dist.

L^2 can't select features L' is a feature selector.

L' is robust to outliers (sq term in L^2 vs mod in L').

L^2 is comp. efficient (L' is not or non-sparse (any))

Norm penalties as constrained optimization.

$$\tilde{J} = J + \alpha \cdot \|\omega\|$$

Equivalent to $\begin{aligned} & \min J \\ \text{s.t. } & \alpha \cdot \|\omega\| \leq k. \end{aligned}$

\therefore Lagrangian $L = J + \alpha (\|\omega\| - k)$

Equivalent to norm-penalties.

if $\|\cdot\|$ is $L^2 \rightarrow$ weights are constrained to be inside the L^2 -norm ball.

Similarly for L^1 .

One can put norm-constraint on each column of ω matrix - different layers have different weights.

Strategy 2: Noise injection.

more data \Rightarrow lower bias & variance

Dataset augmentation: Create fake data.

Example of digit recognition: Transform
to create more points.

Be careful not to mix 6 & 9., b & d.

If we know $P(x_i, y)$ can generate - Egg &
chicken problem.

Noise injection can also be done at the
hidden-unit level.

Noise
Dropout - multiplication by noise.

Adding noise to weights is also possible

stochastic implementation of Bayesian inference
over weights.

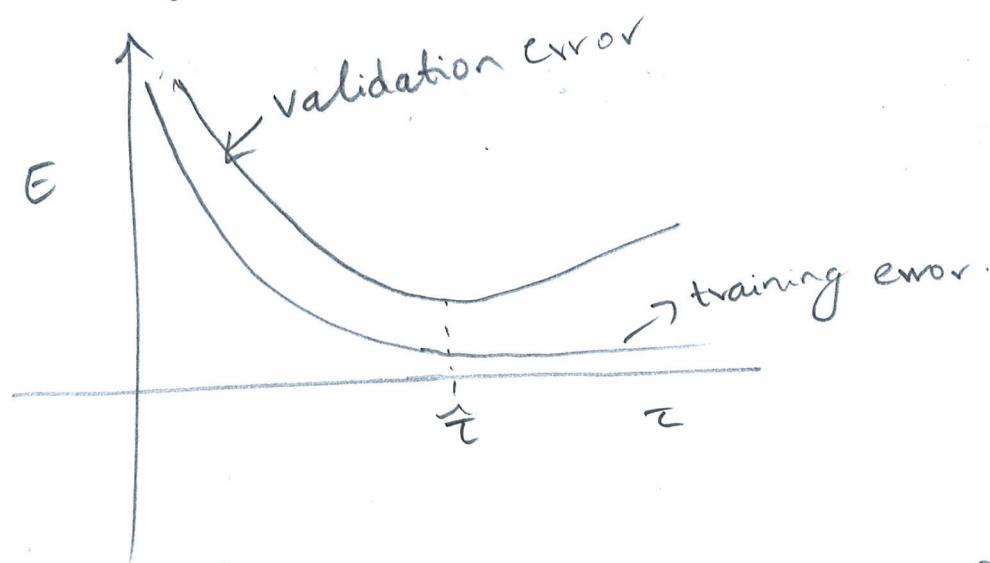
Adding noise is also equivalent to regularization.

Label-smoothing - injecting noise at the output
if the y is correct with $\frac{1}{K}$ + ϵ , Eg: softmax
(18). instead of $\frac{e^{y_i}}{\sum e^{y_i}}$

Regularization - Continued.

Early stopping : Training involves iterative reduction of cost.

Typical training curve:



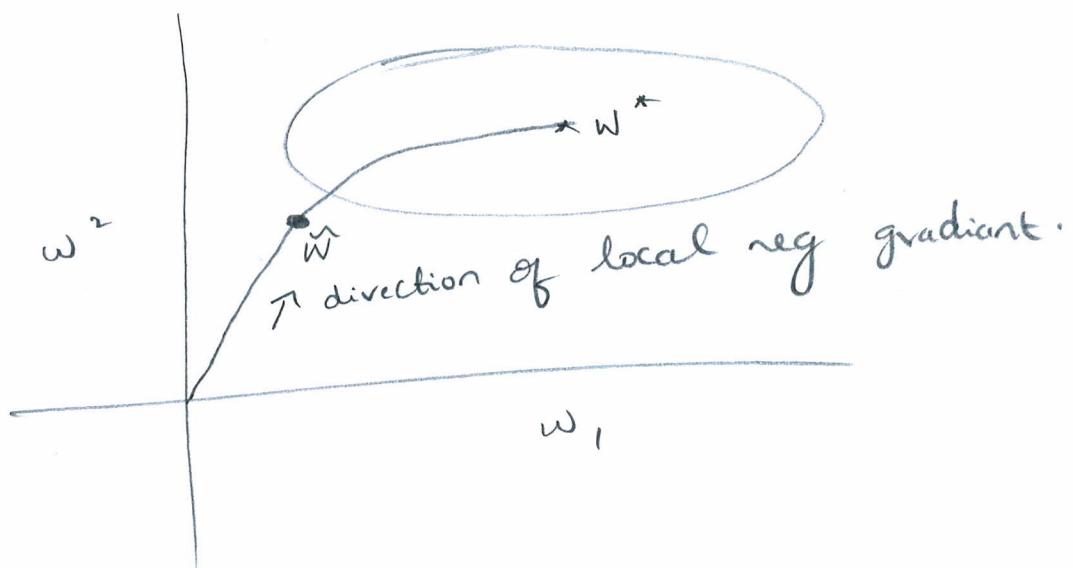
Validation error increases beyond a point \therefore of overfitting.

Early stopping: stop, as VE starts increasing unobtrusive way of regularization (requires no change in obj. fn).

Can be used with other regularizer as well.
Can take this as init value & retrain the whole model with all the training data.

(1)

Why is it a regularizer?



E.S acts as a weight decay. In fact the precise relationship is derived as follows.
Let ω^τ be the weight-update at τ^{th} iteration.

$$\omega^\tau = \omega^{\tau-1} - \epsilon \nabla_{\omega} \hat{f}(\omega^{\tau-1}). \quad (1)$$

$$\hat{f} = f(\omega^*) + \frac{1}{2} (\omega - \omega^*)^T H (\omega - \omega^*).$$

$$(1) = \omega^{\tau-1} - \epsilon H (\omega^{\tau-1} - \omega^*) \quad (2).$$

$$\therefore \nabla_{\omega} \hat{f} = H (\omega - \omega^*).$$

$$\omega^\tau - \omega^* = (I - \epsilon H) (\omega^{\tau-1} - \omega^*)$$

(2).

1/2

$$H = Q \Lambda Q^T \quad [\text{EVD}]$$

$$\therefore \omega^\tau - \omega^* = (I - \epsilon Q \Lambda Q^T)(\omega^{\tau-1} - \omega^*)$$

$$Q^T(\omega^\tau - \omega^*) = (I - \epsilon \Lambda) Q^T(\omega^{\tau-1} - \omega^*) - (1)$$

Now, Let $\omega^0 = 0$.

$$\text{Then, } Q^T \omega^\tau = [I - (I - \epsilon \Lambda)^\tau] Q^T \omega^* - (1)$$

$$\hat{\omega} = Q(\Lambda + \alpha I)^{-1} \Lambda Q^T \omega^*$$

Comparing,

$$Q^T \hat{\omega}_{L_2} = (\Lambda + \alpha I)^{-1} \Lambda Q^T \omega^*$$

$$Q^T \hat{\omega}_{L_2} = [(I - \Lambda + \alpha I)^{-1} \alpha] Q^T \omega^* - (2)$$

\Rightarrow If $\epsilon, \alpha \in \mathbb{C}$ are chosen s.t.

Matrix cookbook.

p6-17.

$$(I - \epsilon \Lambda)^\tau = (\Lambda + \alpha I)^{-1} \alpha.$$

then L_2 & \mathcal{E}_S are equivalent

(3).

Taking log & using Taylor Series for $\log(1+x)$,
 if α_i are small s.t. ($\epsilon \alpha_i \ll 1$ & $\alpha_i/\alpha \ll 1$) Then

$$\tau \approx \frac{1}{\epsilon} \alpha$$

$$\text{or } \alpha \approx \frac{1}{\tau \epsilon}$$

[not clear yet] $\tau \log(1 - \epsilon \alpha) = -1 \log(1 + \alpha) + \log \alpha.$

$$\tau \log(1 - \epsilon \alpha) = -1 \log(1 + \alpha) + \log \alpha.$$

$$\tau + \frac{-1(\epsilon \alpha)}{1 - \epsilon \alpha} = -1. \quad]$$

$\epsilon \rightarrow$ learning rate, $\tau \rightarrow$ iteration index.

$\tau \epsilon$ behaves as effective capacity.

Kind of reciprocal of weight decay
 [inverse of capacity].

So far, we have seen two approaches for regularization: Parameter norm-penalties
Controlling the effective capacity

Third important approach: Confusion is good!!

Noise induction: Go to 76 for intro
after Bagging.

Another approach for regularization: Democracy.
Netflix grand prize
Committees of networks or bagging: predict movie ratings.

Demo with the example of 9, 6 vs 8.

Bagging does better than individual classifiers

Let $h(x)$ be the tree fn we are seeking
to approximate.

Let them be L networks learning this
function.

Let $y_i(x) = h(x) + \epsilon_i(x)$ be the o/p of each network where $\epsilon_i(x)$ is some error fn.

Now, Avg sum-of-sq. error for $y_i(x)$

$$\text{if } E_i = \mathbb{E} \left[\{y_i(x) - h(x)\}^2 \right] = \mathbb{E} [\epsilon_i^2]$$

$$\mathbb{E} [\epsilon_i] = \int \epsilon_i^2(x) p(x) dx$$

Now avg error made by all acting individually

$$\text{if } E_{AV} = \frac{1}{L} \sum_{i=1}^L E_i = \frac{1}{L} \sum_{i=1}^L \mathbb{E} [\epsilon_i^2]$$

Lets now introduce a simple bagging scheme.

$$y_{com}(x) = \frac{1}{L} \sum_{i=1}^L y_i(x)$$

avg of all o/p's.

Now, error due to the committee is

$$\bar{E}_{COM} = \mathbb{E} \left[\left(\frac{1}{L} \sum_{i=1}^L y_i(x) - h(x) \right)^2 \right]$$

$$= \mathbb{E} \left[\left(\frac{1}{L} \sum_{i=1}^L \epsilon_i \right)^2 \right]$$

(6)

^{1/2} Assume $E[\epsilon_i(x)] = 0 \quad E[\epsilon_i \epsilon_j] = 0$

$$\text{Then } E_{\text{com}} = \frac{1}{L^2} \sum_{i=1}^L E[\epsilon_i^2] = \frac{1}{L} E_{\text{av}}$$

\Rightarrow there is a dramatic reduction in error by a factor of L^2 .

not possible in practice $\because E[\epsilon_i \epsilon_j] \neq 0$.

However, Cauchy's inequality

$$\Rightarrow \left(\sum_{i=1}^L \epsilon_i \right)^2 \leq L \sum_{i=1}^L \epsilon_i^2$$

$$\Rightarrow E_{\text{com}} \leq E_{\text{av}}$$

$$\text{One can also take } Y_{\text{com}} = \sum_{i=1}^L d_i y_i(x)$$

for which it can be shown

$$E_{\text{av}} \leq E_{\text{com}}$$

Problem: Too many networks to be trained
-memory!

Another problem: redundancy.

84

In general if $E[\epsilon_i^2] = v$ & $E[\epsilon_i \epsilon_j] = c$

Then avg error = $\frac{1}{K} \sum_i \epsilon_i$

Expected sq. error

$$\begin{aligned} E\left[\frac{1}{K} \sum_i \epsilon_i\right]^2 &= \frac{1}{K^2} E\left[\sum_i \left(\epsilon_i^2 + \sum_{j \neq i} \epsilon_i \epsilon_j\right)\right] \\ &= \frac{1}{K} v + \frac{K-1}{K} c. \end{aligned}$$

if the classifiers are perfectly correlated
 ~~$c=v$~~ , then model avg doesn't help.

perfectly un-correlated case, error = $\frac{v}{K}$

\Rightarrow error $\propto \frac{1}{K}$

Bagging requires datasets sampled with replacement.

(8).

Parameter sharing: Other way of expressing our prior knowledge.

(Creating dependencies b/w model parameters)

Now, we showed, $L_2 \equiv \text{MAP}$ with Gaussian prior

$$P(\omega) = \frac{1}{(2\pi)^{1/2}} \exp\left(-\frac{\omega^2}{2}\right)$$

$$L = \prod_{i=1}^n P(\omega_i) = \frac{1}{(2\pi)^{n/2}} \exp\left\{-\frac{1}{2} \sum_i \omega_i^2\right\}$$

$\Rightarrow L_2$ is encouraging weights to be inside a circle around origin.

Now, instead of one group, ^{diff} weights can lie in different groups.

$$\text{eg: } P(\omega) = \prod_{j=1}^m d_j \phi_j(\omega) \quad [\text{G MAP prior}]$$

\rightarrow known to have better reg effects.

Another example:

Let there be two models A & B with parameters $\omega^{(A)}$ & $\omega^{(B)}$.

$$\hat{y}^{(A)} = f(\omega^{(A)}, x) \quad \hat{y}^{(B)} = g(\omega^{(B)}, x)$$

86 Two different tasks with same features.
If the tasks are similarly, $w^{(A)}$ should be close to $w^{(B)}$

New regularizer $\rightarrow (w^A, w^B)$.

$$= \|w^A - w^B\|_2^2.$$

One approach: Train one classifier under a supervised paradigm by close to another
One trained as un-supervised.

Another powerful approach to Regularization:

Confusion is a good thing:

Noise injection: Go to page 76.

Suppose ridge are perturbed by with a random noise \rightarrow can show that it is similar to L₂ regularizer with L₂ defined by noise covariance.

The Pinnacle of all regularizers:

The Dropout:

We have seen that Bagging always decreases the error.

However, with DNNs, averaging with many separately trained DNNs is prohibitively expensive.

Another problem with Bagging: One has to train with different subsets of data
 But training many with different subsets is not possible \because of limited data.

Further, similar networks should share parameters.

Dropout addresses all of them.

(ii)

Dropout - learning a collection of 2^n possible subnetworks in a DNN
for each perturbation of a training car, a new
thinned car is sampled & trained.

There are 2^n possible subnetworks in a DNN

Simplest car: can be extended with a
probability p.

if off nodes
with (dropping training) will be
Dropout: Arbitrarily remove the hidden

Motivation:
Testing: If a unit is activated with $p = \frac{1}{4}$, it weights are multiple
of their original trained very easily, if at all.
weight sharing, where each trained
thinned network with different
dropout often trained very easily, if at all.
This makes the expected LP (under diff. for dropout).

motivation for dropout: Sexual reproduction

\downarrow

$\frac{1}{2}$ from each parent
(still robust)

criterion for natural selection \neq individual fitness
but miscability of genes.

Ability of genes to work ^{well} with random set of genes makes them robust.

Another motivation: Ten successful conspiracies involving 5 people each is a better way to create havoc than one big conspiracy that begins so ppl all playing their part correctly!

Formalization: Multiplication of hidden unit with Bernoulli noise.

X

(15).

90

$$X \in \mathbb{R}^{N \times D}, y \in \mathbb{R}^N, w \in \mathbb{R}^D$$

$$\|y - Xw\|^2$$

under dropout,

X is retained with p ,

which can be expressed as

$$R \odot X \quad \text{when } R \in \{0, 1\}^{N \times D}$$

with $R_{ij} \sim \text{Bernoulli}(p)$.

$$\therefore \text{Obj fn: } \min_w E_{\text{Bernoulli}(p)} [\|y - (R \odot X)w\|^2]$$

$$\Rightarrow \min_w \|y - Rxw\|^2 + p(1-p) \|Tw\|^2$$

$$T = (\text{diag}(x^T x))^{1/2}$$

Dropout \equiv ridge with $T = I$

T scales the weight loss with std dev of data.

if a ~~row~~ data ~~dim~~ var. a lot,
it will ~~square~~ the weight now

Another formulation.

$$\min_w \|y - x^T w\|^2 + \frac{1-p}{p} \|Tw\|^2$$

$\tilde{w} = pw$, for $p \approx 1$, all ips are retained
 & with more dropout ($p < c$)
 regularizes to zero

One can also have multiplicative Gaussian noise.

$N(1, 1)$ works better than Bernoulli e.g.

$h_i \sim h_i + h_i' \sim N(0, 1)$

or $h_i' \sim \underline{N}(1, 1)$.

in general $\mathbf{r}' \sim N(1, \sigma^2)$ where σ is another hyper-parameter.

Adv of null-noise: ReLU can't overlook noise.

Power of dropout: adaptive destruction of info
 is occurring at hidden units (higher level)
 of abstractions:

Example of mult vs non.
 Enhancing the model capacity!!!

