

Introduction to Bias & Variance

Let us consider the problem of fitting a curve through a given set of points.

We can consider two models:

Simple $y = \hat{f}(x) = w_1 x + w_0$
(degree: 1)

Complex $y = \hat{f}(x) = \sum_{i=1}^{25} w_i x^i + w_0$
(degree: 25)

Note: In both cases we are making an assumption about how y is related to x . We have no idea about the true relation $f(x)$

We are given 500 data points

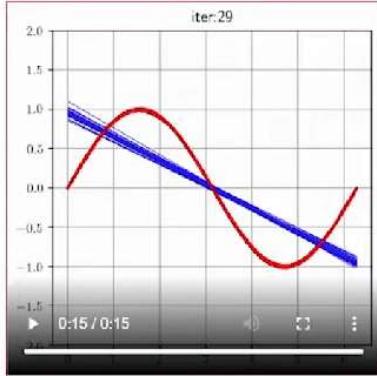
Let's sample 200 data points & train a simple & complex model

We'll repeat the process 'k' times to train multiple models
(each model will see different sample of the data)

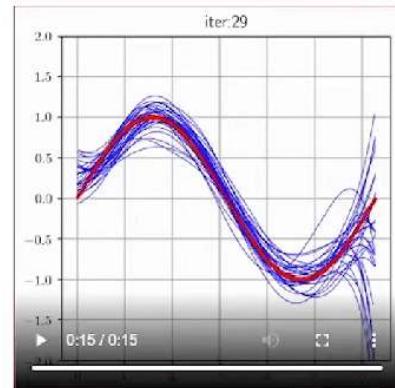
Simple models trained on different samples of data do not differ much from each other.

However they are very far from the true curve (underfitting)

On the other hand, complex models trained on different samples are very different from each other (high variance)

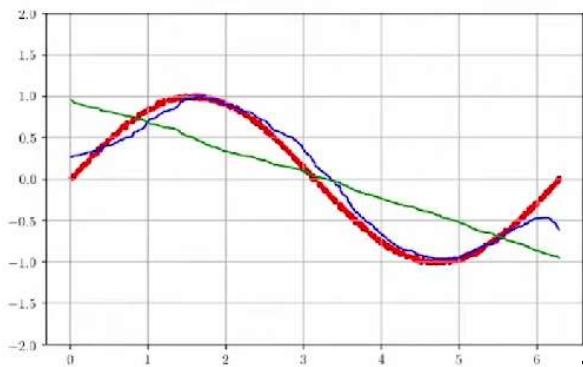


Simple Model



Complex Model

Let $f(x)$ be the true model & $\hat{f}(x)$ be our estimate of the model (simple or complex) then,



Green: Avg values of $\hat{f}(x)$ for simple model

Blue: Avg values of $\hat{f}(x)$ for complex model

Red: True model $f(x)$

$$\text{Bias } \hat{f}(x) = E[\hat{f}(x)] - f(x)$$

(difference b/w avg & true value)

	Simple Model	Complex Model
Bias	Very High	Very Low

$$\text{Variance } (\hat{f}(x)) = E[(\hat{f}(x) - E[\hat{f}(x)])]^2$$

(How much the different $\hat{f}(x)$'s differ from each other)

	Simple Model	Complex Model
Variance	Very Low	Very High

Training Errors vs Testing Errors

Consider a new point (x, y) which was not seen before during training.

If we use the model $\hat{f}(x)$ to predict the value of y then the mean square error is given by:

$$E[(y - \hat{f}(x))^2]$$

We can show that $E[(y - \hat{f}(x))^2] = \text{Bias}^2$

+ Variance

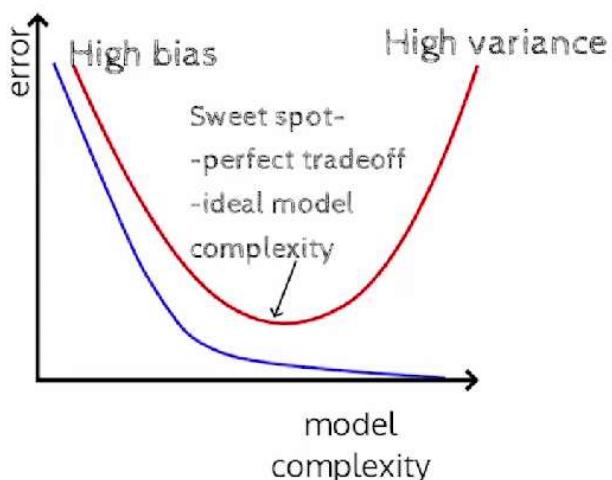
+ σ^2 (irreducible error)

The parameters of $\hat{f}(x)$ (all the w_i 's) are trained using a training set $\{(x_i, y_i)\}_{i=1}^n$.

At test time, we are interested in evaluating the model on validation (unseen) set which was not used for training.

Now we have two entities:

$\text{train}_{\text{error}}$ & $\text{test}_{\text{error}}$



Given: n training points,
 m testing points.

$$\text{train}_{\text{err}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

$$\text{test}_{\text{err}} = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{f}(x_i))^2$$

Estimate Errors from Test Data

Let $D = \{x_i, y_i\}_{i=1}^n$

We know that some true function exists s.t.

$$y = f(x) + \varepsilon$$

means that y_i is related to x_i by some function f
but there is some noise ε in the relation

Let's assume $\varepsilon \sim N(0, \sigma^2)$

We will use \hat{f} to approximate f & estimate the parameters
of training data TCD.

$$y = \hat{f}(x)$$

We are interested in knowing

$$E[(\hat{f}(x) - f(x))^2]$$

but we cannot estimate this directly as we do not
know f

$$E[(\hat{y} - y)^2] = E[(\hat{f}(x) - f(x) - \varepsilon)^2]$$

$$= E[(\hat{f}(x) - f(x))^2 - 2\varepsilon(\hat{f}(x) - f(x)) + \varepsilon^2]$$

$$= E[(\hat{f}(x) - f(x))^2] - 2E[\varepsilon(\hat{f}(x) - f(x))] + E[\varepsilon^2]$$

$$E[(\hat{f}(x) - f(x))^2] = E[(\hat{y} - y)^2] - E[\varepsilon^2] + 2E[\varepsilon(\hat{f}(x) - f(x))]$$

Case 1: Using Test Data

$$- E[(\hat{y}_i - y_i)^2] = \frac{1}{n+m} \sum_{i=1}^{n+m} (\hat{y}_i - y_i)^2 - \sigma^2 + 2E[\varepsilon(\hat{f}(x) - f(x))]$$

$$\underbrace{E[(\hat{f}(x) - f(x))^2]}_{\text{true error}} = \underbrace{\frac{1}{m} \sum_{i=n+1}^{n+m} (\hat{y}_i - y_i)^2}_{\text{empirical estimation of error}} - \underbrace{\sigma^2}_{\text{small constant}} + 2E[\epsilon(\hat{f}(x) - f(x))]$$

covariance (ϵ , $(\hat{f}(x) - f(x))$)

$$\begin{aligned}\text{cov}(X, Y) &= E[(X - \mu_X)(Y - \mu_Y)] \\ &= E[X(Y - \mu_Y)] \quad [\text{if } \mu_X = E[X] = 0] \\ &= E[XY] - E[X]\mu_Y \\ &= E[XY] - \mu_Y E[X] \\ &= E[XY]\end{aligned}$$

$$\text{Now } \epsilon = y - f(x)$$

This $y - f(x)$ is independent of $(\hat{f}(x) - f(x))$
as it is based on test data & the latter is train data.
∴ y is independent of $\hat{f}(x)$
∴ ϵ is independent of $(\hat{f}(x) - f(x))$
∴ covariance will be 0

∴ true error = empirical test error + small constant

Case 2 : Using Training Data

$$\underbrace{E[(\hat{f}(x) - f(x))^2]}_{\text{true error}} = \underbrace{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2}_{\text{empirical estimation of error}} - \underbrace{\sigma^2}_{\text{small constant}} + 2E[\epsilon(\hat{f}(x) - f(x))]$$

covariance (ϵ , $(\hat{f}(x) - f(x))$)

The error

estimation of
error

$$E[\epsilon \cdot (\hat{f}(x) - f(x))] \neq E[\epsilon] \cdot E[(\hat{f}(x) - f(x))] \neq 0$$

\therefore empirical train error is smaller & does not give a true picture

True Error vs Model Complexity

According to Stein's lemma we can show that,

$$\frac{1}{n} \sum_{i=1}^n \varepsilon_i (\hat{f}(x_i) - f(x_i)) = \frac{\sigma^2}{n} \sum_{i=1}^n \frac{\partial \hat{f}(x_i)}{\partial y_i}$$

When will $\frac{\partial \hat{f}(x_i)}{\partial y_i}$ be high?

When a small change in the observation causes a large change in the estimation (\hat{f})

A complex model will be more sensitive to changes in observations whereas a simple model will be less sensitive

Hence we can say:

$$\text{true error} = \underset{\text{train error}}{\text{empirical}} + \underset{\text{constant}}{\text{small}} + \underset{\substack{\leftarrow \text{model complexity}}}{\Omega}$$

Hence while training, instead of minimizing $L_{\text{train}}(\theta)$ we should minimize :

$$\min_{\text{wrt } \theta} L_{\text{train}}(\theta) + \Omega(\theta) = L(\theta)$$

$\Omega(\theta)$ acts as an approximate for $\frac{\sigma^2}{n} \sum_{i=1}^n \frac{\partial \hat{f}(x_i)}{\partial y_i}$

This is the basis of regularization methods.

This is the basis of regularization methods.

L2 Regularization

Regularization term \rightarrow Proxy for model complexity.

For L2 regularization we have:

$$\bar{L}(\omega) = L(\omega) + \frac{\alpha}{2} \|\omega\|^2$$

$$\text{For SGD : } \nabla \bar{L}(\omega) = \nabla L(\omega) + \alpha \omega$$

$$\text{Update rule: } \omega_{t+1} = \omega_t - \eta \nabla L(\omega_t) - \eta \alpha \omega_t$$

Assume ω^* is the optimal solution for $L(\omega)$ (not $\bar{L}(\omega)$)
i.e. solution in the absence of regularization (ω^* optimal $\rightarrow \nabla L(\omega^*) = 0$)

$$\text{Consider } \omega = \omega^* + u$$

$$\therefore u = \omega - \omega^*$$

$$L(\omega^* + u) = L(\omega^*) + u^T \nabla L(\omega^*) + \frac{1}{2} u^T H u$$

$$L(\omega) = L(\omega^*) + (\omega - \omega^*)^T \nabla L(\omega^*) + \frac{1}{2} (\omega - \omega^*)^T H (\omega - \omega^*)$$

$$= L(\omega^*) + \frac{1}{2} (\omega - \omega^*)^T H (\omega - \omega^*) \quad [\nabla L(\omega^*) = 0]$$

$$\nabla L(\omega) = \nabla L(\omega^*) + H(\omega - \omega^*)$$

[H \rightarrow Hessian Matrix]

$$= H(\omega - \omega^*)$$

Now,

$$\nabla \bar{L}(\omega) = \nabla L(\omega) + \alpha \omega$$

Let $\bar{\omega}$ be the optimal solution for $\bar{L}(\omega)$

$$\therefore \nabla \bar{L}(\bar{\omega}) = 0$$

Let w^* be the optimal solution for $\min_{w \in \mathbb{R}^n}$

$$\therefore \nabla L(\bar{w}) = 0$$

$$H(\bar{w} - w^*) + \alpha \bar{w} = 0$$

$$\therefore (H + \alpha I) \bar{w} = H w^*$$

↑
identity
vector

$$\therefore \bar{w} = (H + \alpha I)^{-1} H w^*$$

as $\alpha \rightarrow 0$, $\bar{w} \rightarrow w^*$ [no regularization]

Let us analyze when $\alpha \neq 0$

If H is a symmetric Positive Semi Definite Matrix,

$$H = Q \Lambda Q^T \quad [Q \text{ is orthogonal, } Q Q^T = Q^T Q = I]$$

$$\begin{aligned}\therefore \bar{w} &= (H + \alpha I)^{-1} H w^* \\ &= (Q \Lambda Q^T + \alpha I)^{-1} Q \Lambda Q^T w^* \\ &= (Q \Lambda Q^T + \alpha Q \Lambda Q^T)^{-1} Q \Lambda Q^T w^* \\ &= [Q(\Lambda + \alpha I) Q^T]^{-1} Q \Lambda Q^T w^* \\ &= Q^{-1} (\Lambda + \alpha I)^{-1} Q^T Q \Lambda Q^T w^* \\ &= Q(\Lambda + \alpha I)^{-1} \Lambda Q^T w^* \quad [Q^{-1} = Q]\end{aligned}$$

when $\alpha = 0$, w^* becomes the optimal solⁿ as

$$\therefore \bar{w} = (Q D Q^T) w^* \quad [D = (\Lambda + \alpha I)^{-1} \Lambda, \text{ a diagonal matrix}]$$

w^* acts rotated first by Q^T to give $Q^T w^*$ $[Q \rightarrow n \times n]$

w* gets rotated first by Q^T to give $Q^T w^*$
 If $\alpha=0$, Q rotates $Q^T w^*$ back to give w* matrix

$$\begin{cases} Q \rightarrow n \times n \\ w^* \rightarrow n \times 1 \\ D \rightarrow n \times n \end{cases}$$

If $\alpha \neq 0$, $(\Lambda + \alpha I)^{-1} = \begin{bmatrix} \frac{1}{\lambda_1 + \alpha} & & & \\ & \frac{1}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{1}{\lambda_n + \alpha} \end{bmatrix}$ Each element i of $Q^T w^*$ gets scaled by $\frac{\lambda_i}{\lambda_i + \alpha}$ before it is rotated back by Q

$$D = (\Lambda + \alpha I)^{-1} = \begin{bmatrix} \frac{\lambda_1}{\lambda_1 + \alpha} & & & \\ & \frac{\lambda_2}{\lambda_2 + \alpha} & & \\ & & \ddots & \\ & & & \frac{\lambda_n}{\lambda_n + \alpha} \end{bmatrix}$$

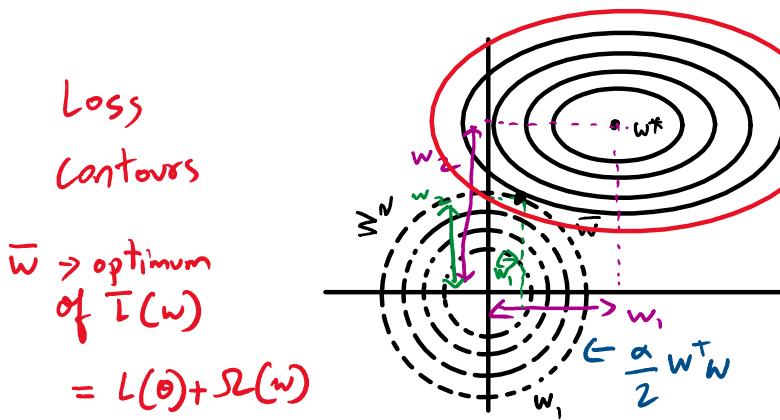
If $\lambda_i \gg \alpha$ then, $\frac{\lambda_i}{\lambda_i + \alpha} = 1$

If $\lambda_i \ll \alpha$ then, $\frac{\lambda_i}{\lambda_i + \alpha} = 0$

\therefore only significant directions will be retained,

Effective Parameters = $\sum_{i=1}^n \frac{\lambda_i}{\lambda_i + \alpha} < n$

& complexity has reduced!



The weights have shrunk in presence of regularization.

w_1 has decreased more than w_2 coz it is more imp than w_2

$$= L(\theta) + \frac{\alpha}{2} w^T w$$

Even the loss we want to be small.

If w_2 shrank loss would be higher \rightarrow see red contours

Dataset Augmentation

Given an image of handwritten digit, we can rotate it by n degrees k times, shift positions, blur it & also change pixels. This will still not change the label.

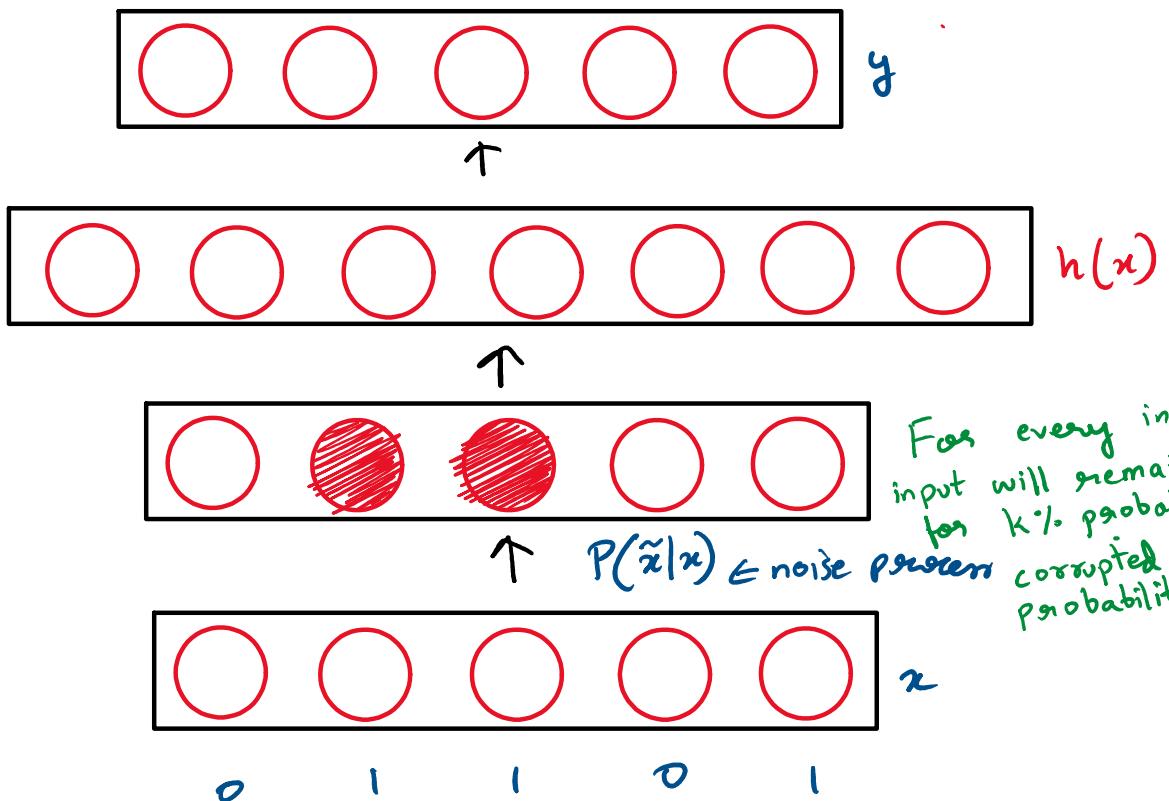
∴ We exploit the fact that transformations to the image do not change the image & create more training data. This is called data augmentation. Data created using knowledge of the task.

Parameter Sharing

Used in CNNs.

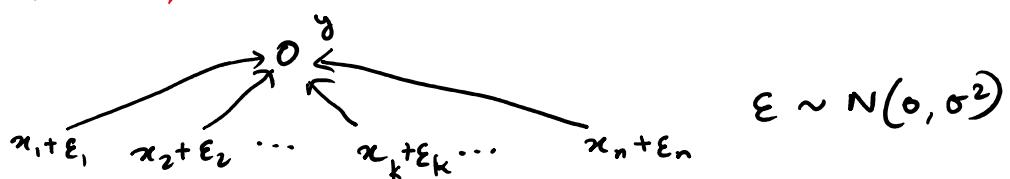
Same filter is applied to different positions of the image.
Or same weight matrix acts on different input neurons.

Injecting Noise At Inputs



For every epoch, the corruption happens differently.

$$\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n$$



$$\begin{aligned}
 \tilde{x}_i &= x_i + \epsilon_i & \tilde{y} &= \sum_{i=1}^n w_i \tilde{x}_i \\
 & & &= \sum_{i=1}^n w_i x_i + \sum_{i=1}^n w_i \epsilon_i \\
 & & \boxed{\tilde{y} = \hat{y} + \sum_{i=1}^n w_i \epsilon_i}
 \end{aligned}$$

$$\text{Now, } E[(\tilde{y} - y)^2] = E\left[\left(\hat{y} + \sum_{i=1}^n w_i \epsilon_i - y\right)^2\right]$$

$$\begin{aligned}
 \text{Now, } E[(\tilde{y} - y)^2] &= E\left[\left(\hat{y} + \sum_{i=1}^n w_i \epsilon_i - y\right)^2\right] \\
 &= E\left[(\hat{y} - y)^2 + \left(\sum_{i=1}^n w_i \epsilon_i\right)^2\right] \\
 &= E[(\hat{y} - y)^2] + E\left[\left(\sum_{i=1}^n w_i \epsilon_i\right)^2\right] \\
 &= E[(\hat{y} - y)^2] + E\left[2(\hat{y} - y) \sum_{i=1}^n w_i \epsilon_i\right] + E\left[\left(\sum_{i=1}^n w_i \epsilon_i\right)^2\right] \\
 &= E[(\hat{y} - y)^2] + 0 + E\left[\sum_{i=1}^n w_i^2 \epsilon_i^2\right] \quad \boxed{\begin{array}{l} \epsilon_i \text{ is independent of} \\ \epsilon_j \& \epsilon_i \text{ is} \\ \text{independent of } (\hat{y} - y) \end{array}} \\
 &= E[(\hat{y} - y)^2] + \sigma^2 \sum_{i=1}^n w_i^2 \quad \boxed{E[\epsilon_i^2] = \sigma^2}
 \end{aligned}$$

∴ In a simple output network, w/o any non-linearity, adding noise to the inputs is the same as using weight decay!
 aka L2 norm penalty.

Adding Noise to the Output

For a given classification:



y	0	1	2	3	4	5	6	7	8	9
	0	0	1	0	0	0	0	0	0	0

Hard Targets

If we try to minimize the cross entropy loss:

$$\text{minimize: } \sum_{i=0}^q p_i \log q_i \quad p \rightarrow \text{true label} \\ q \rightarrow \text{predicted label}$$

true distribution $\rightarrow p: 0010000000$

estimated distribution $\rightarrow q$

Intuition

Don't trust the label as it is noisy

Instead, add soft targets



Taking a small mass from probability & adding it to all other labels & they'll be my soft targets

$\frac{\varepsilon}{9}$	$\frac{\varepsilon}{9}$	$1 - \varepsilon$	$\frac{\varepsilon}{9}$						
-------------------------	-------------------------	-------------------	-------------------------	-------------------------	-------------------------	-------------------------	-------------------------	-------------------------	-------------------------

Soft Targets

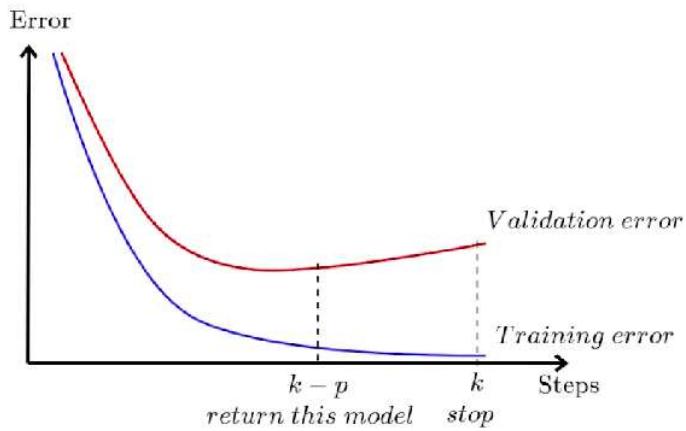
$\varepsilon \rightarrow$ small positive constant

$$\text{minimize: } \sum_{i=1}^q p_i \log q_i$$

true distribution + noise. $0 = \sum \varepsilon \cdot (1 - \varepsilon) \dots \varepsilon \cdot \varepsilon$

True distribution + noise: $P = \left\{ \frac{\epsilon}{q}, \frac{\epsilon}{q}, 1-\epsilon, \dots, \frac{\epsilon}{q} \right\}$
estimated distribution: q

Early Stopping



We want to find the sweet spot where training error is low & validation error is also as low as possible.

We start by tracking the validation error.

Have a patience parameter p

After p consecutive steps, if validation error doesn't decrease & continues to increase (remains monotonic) then we'll have to stop training.

Suppose we are at step k & for the last p steps, validation error doesn't decrease, we'll throw away the last p updates & retain the model at $(k-p)$ steps!

This is one of the most widely used form of regularization. Can also be used w/ other regularizers (ℓ_2)

How does this act as a regularizer?

$$\begin{aligned} \text{The update rule of SGD: } w_{t+1} &= w_t - \eta \nabla w_t \\ &= w_0 - \eta \sum_{i=1}^t \nabla w_i \end{aligned}$$

Let c be the maximum value of ∇w_i then,

Let z be the maximum value of ∇w_i then,

$$|w_{t+1} - w_0| \leq \eta t |z|$$

t is controlling how far w_t can go from w_0

According to Taylor Series:

$$\begin{aligned} L(w) &= L(w^*) + (w-w^*)^T \nabla L(w^*) + \frac{1}{2} (w-w^*)^T H(w-w^*) \\ &= L(w^*) + \frac{1}{2} (w-w^*)^T H(w-w^*) \quad [w^* \text{ is optimal}] \\ &\quad \therefore \nabla L(w^*) = 0 \end{aligned}$$

Taking derivative on both sides:

$$\nabla L(w^*) = H(w-w^*)$$

The SGD update rule is:

$$\begin{aligned} w_{t+1} &= w_t - \eta \nabla w_t \\ &= w_t - \eta H(w_t - w^*) \\ w_{t+1} &= (I - \eta H)w_t + \eta H w^* \end{aligned}$$

Using Eigen Value Decomposition of H as $H = Q \Lambda Q^T$,

$$w_{t+1} = (I - \eta Q \Lambda Q^T)w_t + \eta Q \Lambda Q^T w^*$$

If we start with $w_0 = 0$ then we can show,

$$w_{t+1} = Q [I - (I - \eta \Lambda)^t] Q^T w^*$$

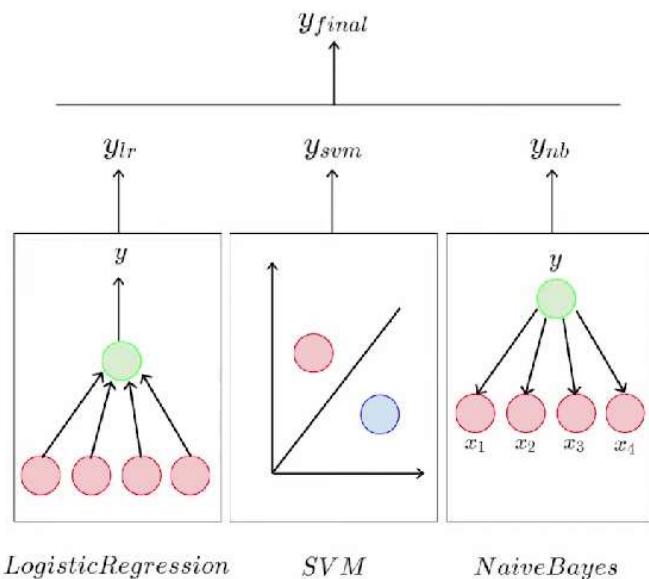
$$w_{t+1} = Q \left[I - (I - \epsilon \Lambda)^T \right] Q^T w^*$$

We observe that $w_t = \tilde{w}$ if we choose ϵ, t & α such that,

$$(\mathbb{I} - \epsilon A)^t = (A + \alpha \mathbb{I})^{-1} \alpha$$

Ensemble Methods

Idea: Combine the output of different models to reduce generalization error. The models can correspond to different classifiers.



It could also be different instances of the same classifier, trained with:
different hyperparameters
different features
different samples of training data
This is called Bagging.

Bagging: form an ensemble using different forms of the same classifier.

From a given dataset, construct multiple training sets by sampling with replacement (T_1, T_2, \dots, T_k)

Training i^{th} instance of the classifier using training set T_i

When would bagging work?

Consider a set of k Logistic Regression models.

Suppose that each model makes an error ϵ_i on a test example.

$$1 - L = 1 - (1 - \epsilon_1)(1 - \epsilon_2) \dots (1 - \epsilon_k) = 1 - (1 - k\epsilon)^k$$

Example .

Let ϵ_i be drawn from a zero mean multivariate normal distribution. $\epsilon_i \sim N(0, \sigma^2)$

$$\text{Variance} = E[\epsilon_i^2] = \sigma^2$$

$$\text{Covariance} = E[E_i, E_j] = C$$

The error made by the average prediction of all the models is $\frac{1}{k} \sum_i e_i$ ← from one sample ← from all samples

The expected savare error: $mse = E\left[\left(\frac{1}{k} \sum_i \varepsilon_i^2\right)^2\right]$

$$= \frac{1}{k^2} E \left[\sum_i \varepsilon_i^2 + \sum_i \sum_{i \neq j} \varepsilon_i \varepsilon_j \right]$$

$$= \frac{1}{k^2} \left(\sum_i E[\varepsilon_i^2] + \sum_i \sum_{i \neq j} E[\varepsilon_i \varepsilon_j] \right)$$

↑
Variance

↑
Covariance

$$= \frac{1}{k^2} (kv + k(k-1)c)$$

$$mse = \frac{1}{k} v + \frac{k-1}{k} c$$

When would bagging work?

If the errors of the model are perfectly correlated, then $V = C$ & $mse = V$ which is the same as the mse

then $V = C$ & $\text{mse} = V$ which is the same as the mse of any of the model.

If the errors of the model are independent or uncorrelated then $C = 0$ & the $\text{mse} = \frac{1}{k} V$

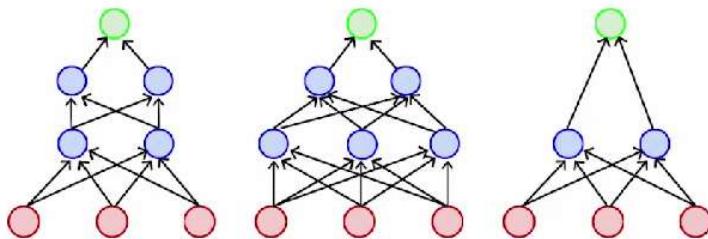
On average, the ensemble will perform at least as well as its individual members

Dropout

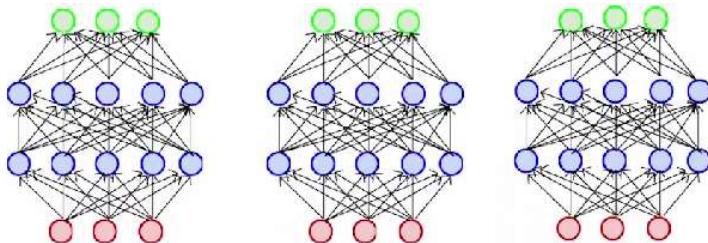
Typically, bagging always helps.

However in the context of deep learning, training several large networks for making an ensemble is expensive.
We have two options:

Option 1: Train several neural networks having different architectures (very expensive)



Option 2: Train multiple instances of the same network, using different training samples (again, expensive)

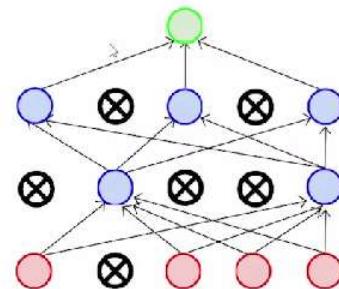
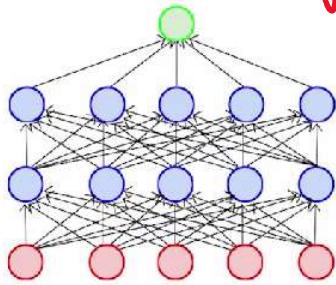


The solution is : Dropout

It refers to dropping out some units (neurons).

Temporarily remove a node & all its incoming & outgoing connections resulting in a thinned network

connections resulting in a thinned network



Each node is retained with a fixed probability ($p=0.5$) for hidden & $p=0.8$ for visible nodes

For a network of n nodes, each node is retained or dropped resulting in a new architecture.

We can construct $\underline{2^n}$ thinned networks from a given neural network

How do we train so many networks?

Trick (1): Share the weights across the networks

Trick (2): Sample a different network for each training instance

Steps

- Initialize all parameters & start training.
- Apply dropout for first training instance (mini batch)
- Compute loss & backpropagate

Which Parameters to update?

- Only the ones who participated !

- Only the ones who participated !

In the second instance (or mini batch), apply dropout again

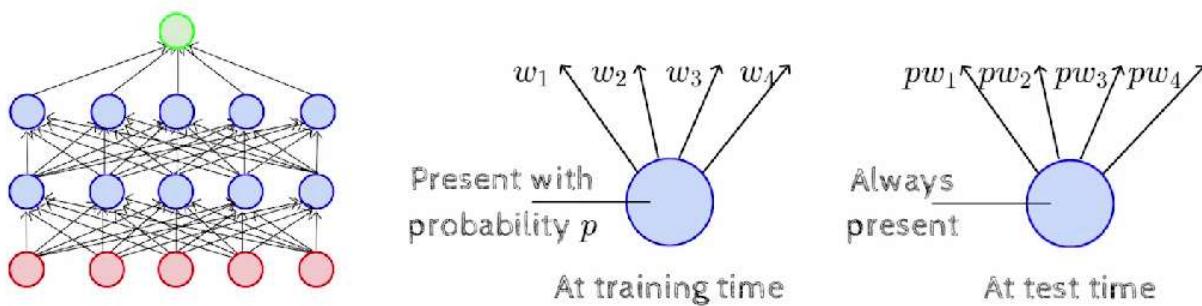
Again compute the loss & backpropagate to the active weights

If a particular weight was active for both training instances then it would be updated twice.

If a particular weight was active for one training instances then it would be updated only once.

What to do during testing ?

If a node was present only with probability p , scale the output by p times.



Dropout, essentially applies a masking noise to the hidden units

It also prevents hidden units from "coadopting"

It also prevents hidden units from "coadopting"
↑
relying on other nodes

Suppose h_i learns to detect a face by firing
on detecting a nose

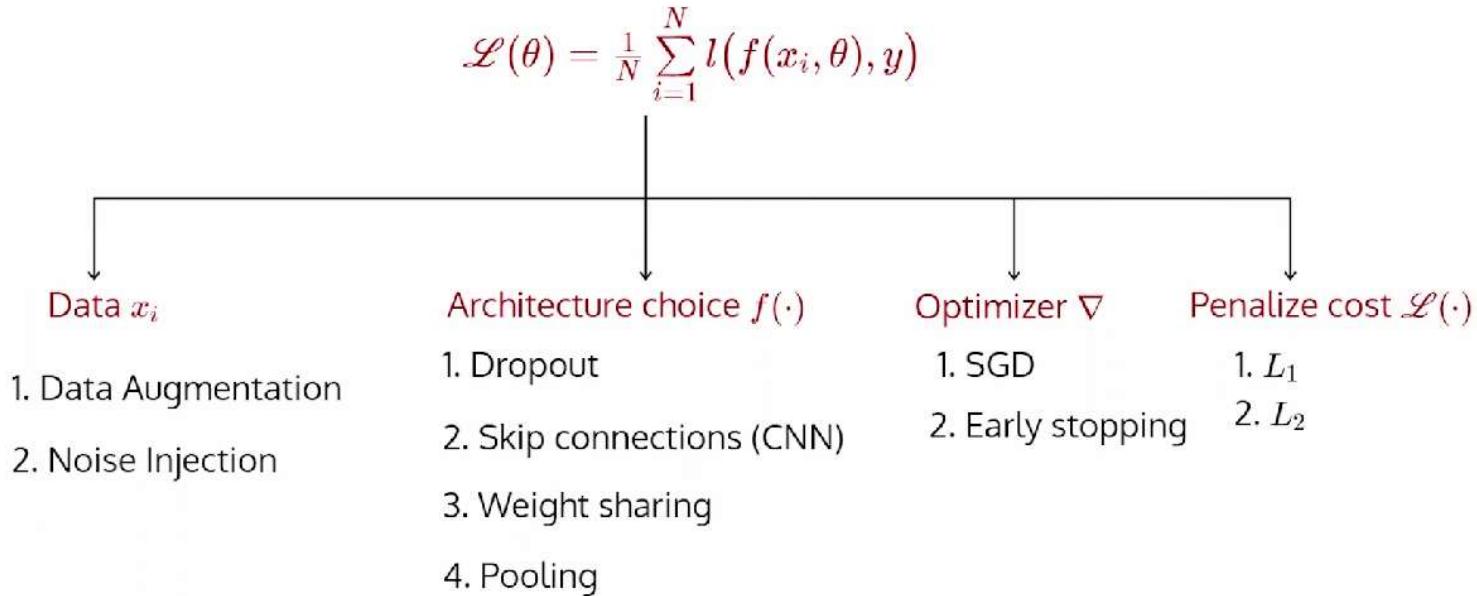
Dropping h_i then corresponds to erasing
the information that a nose exists

The model should then learn another h_i
which redundantly encodes the presence of
a nose

Or the model should learn to detect the
face using other features

Summary

Grouping Regularization Techniques



Grouping Regularization Techniques

