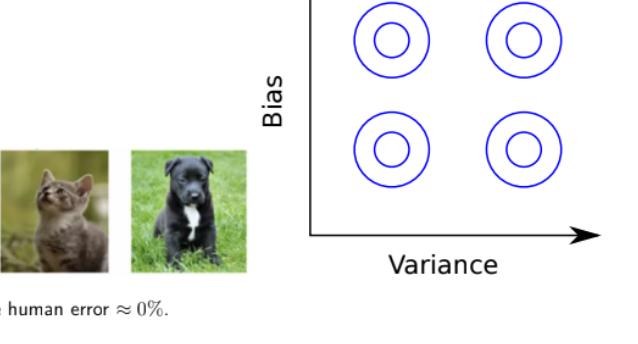
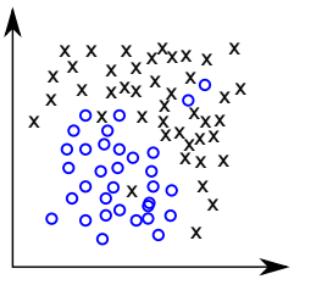


consegue operar nem curso

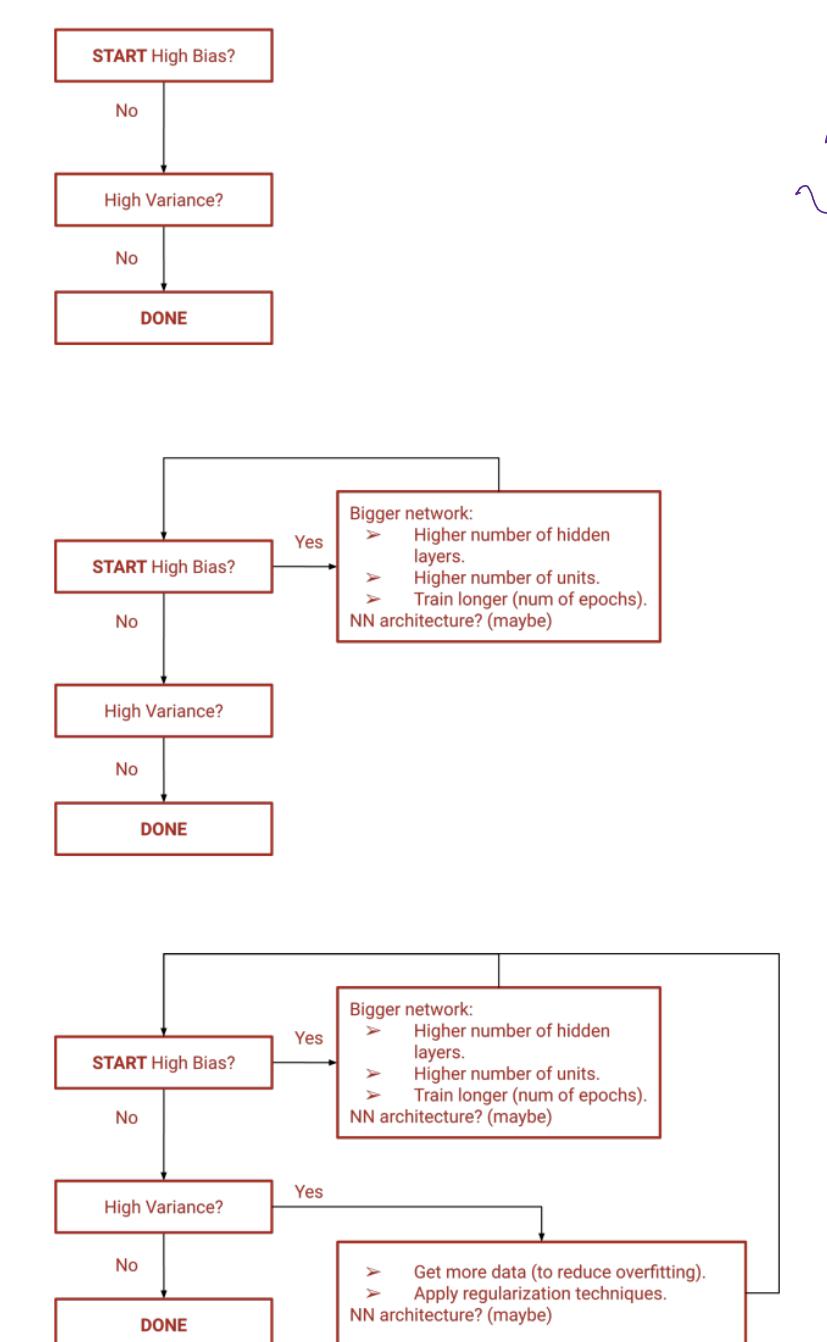


The worst of both worlds (high bias and high variance)



Key Takeaway of Bias-Variance Trade-Off

- A model with greater complexity might be theoretically more accurate (i.e., low bias).
- But we have no control on what it might predict on a new training data set.
- Some of these must be wrong - contradiction of model's accuracy.
- A more accurate model for infinite data is not a more accurate model for finite data.
- No model can be independent of data.



Assume

- x : Independent variable
- y : Dependent variable
- $f(x)$: Describes the true underlying dependence of y on x .
- $\hat{f}(x)$: The result of $f(x)$ and random noise.
- ϵ : Random variable representing noise
- $\hat{y}_i = \hat{f}(x_i) + \epsilon_i$

Recall

$$\begin{aligned} \text{var}(X) &= E[(X - E[X])^2] \\ &= E[X^2] - 2E[X]E[X] + E^2[X] \\ &= E[X^2] - 2E^2[X] + E^2[X] \\ &= E[X^2] - E^2[X] \end{aligned}$$

And also, since $E[\epsilon] = 0$, we have:

$$\begin{aligned} \text{var}(\epsilon) &= E[\epsilon^2] - E^2[\epsilon] \\ \text{var}(\epsilon) &= E[\epsilon^2] - \sigma^2 \end{aligned}$$

Goal

- Model the underlying real-life problem.
- I.e., find f s.t. $\hat{f}(x) \approx f(x)$ by reducing MSE = $E[(Y - f(X))^2]$

Bias

- Difference of average value of prediction mistakes (over different realizations of training data) to the true underlying function $f(x)$ for a given unseen test point.
- i.e. $\text{bias}[\hat{f}(x)] = E[\hat{f}(x)] - f(x)$
- Notice that $f(x)$ is a random variable affected by the randomness in which we obtain training data.

Variance

- Mean Squared Deviation of $\hat{f}(x)$ from its expected value $E[\hat{f}(x)]$ over different realizations of training data, i.e.:
- $\text{var}(\hat{f}(x)) = E[(\hat{f}(x) - E[\hat{f}(x)])^2]$

Summary

- $\text{MSE} = E[(y - \hat{f}(x))^2]$
- $\text{bias}[\hat{f}(x)] = E[\hat{f}(x)] - f(x)$
- $\text{var}(\hat{f}(x)) = E[(\hat{f}(x) - E[\hat{f}(x)])^2]$

* Validation set = portuguese \Rightarrow dataset de treino em 2 conjuntos

↳ 80% 20% etc

↳ Cross-validation

↳ k-fold \Rightarrow dividir em k subsets e pegar o que deixa de k

Um regra, é somar os valores absolutos de todos os pesos

A função de custo regularizada com λ também usa um hiperparâmetro " α "

$$J(w, X, y) \rightarrow \alpha \|w\|^2 + J(w, X, y)$$

\rightarrow podemos notar que tem um fator constante

* Regularization visto Archer

\rightarrow Entrar overfitting \rightarrow diminuir a variação

Bonsamente você adiciona os pesos da custo:

paramétrico regularizado constante

$$\frac{\lambda}{2m} \|w\|_2^2 \rightarrow \text{Norma da quebra de "l1"}$$

$$\|w\|_2^2 = \sum_j w_j^2 = w^T w$$

\rightarrow Regularização \rightarrow menor variação

\rightarrow Você adiciona geralmente peso a w , porque tem mais chances de ter problemas de variação, pelo maior dimensão

\rightarrow λ Reg $\rightarrow \frac{\lambda}{m} \sum_j |w_j| \rightarrow L_1 \text{ norm}$

$\rightarrow w = espaco \rightarrow$ se regras, u vai ter mais zeros

$\rightarrow \lambda \rightarrow$ bairinho \rightarrow Reg. paramétrico baseado no test set que você pode tirar

\rightarrow No entanto, você não quer todos, só quer que o custo function funcione:

$$J(w, b) = \frac{1}{m} \sum_i J(\hat{y}_i, y_i) + \frac{\lambda}{2m} \sum_j \|w_j\|^2$$

onde
 $\|w\|^2 \rightarrow \sum_j w_j^2$ ou seja, $\frac{1}{m} \sum_i \sum_j w_j^2$

\rightarrow o gradient descent vai ficar da seguinte forma:

$$dW^{(t+1)} \rightarrow \text{backprop} + \frac{\lambda}{m} w^{(t)}$$

$d \rightarrow$ update contínuo & normal

\rightarrow Se chama weight decay porque cada vez que for passar a nova dimensão

\rightarrow Se chama weight decay porque cada vez que for passar a nova dimensão

\rightarrow Pode ser weight de ser muito longo, diminuindo ele do peso

sem peso com "longo" e "

Dropout Regularization

Adiciona a probabilidade de novo bairinho ser excluído (dropout), de modo que excluimos

Nov node e os ligados a partir dele, resultando em uma rede menor

\rightarrow Implementação

$\beta = np.random.rand()$ \rightarrow Sendo como moeda, de modo que seu aplicando nos nós,

$\text{if } \beta > 0.5, \text{ então } = 1 \text{ e caso contrário } = 0$

\rightarrow Tamanha se se violar $\beta > 0.5$, então $= 1$ e caso contrário $= 0$

From the Mean Squared Error (MSE):

$$\begin{aligned} \mathbb{E}[(y - f(x))^2] &= \mathbb{E}[(f(x) + \epsilon - f(x))^2] \\ &= \mathbb{E}[(f(x) - f(x))^2] + 2\mathbb{E}[(f(x) - f(x))\epsilon] + \mathbb{E}[\epsilon^2] \\ \text{Since expectation is linear and } \epsilon \text{ and } f \text{ are independent:} \\ &= \mathbb{E}[(f(x) - f(x))^2] + 2\mathbb{E}[(f(x) - f(x))]\mathbb{E}[\epsilon] + \mathbb{E}[\epsilon^2] \\ \text{MSE} &= \mathbb{E}[(y - f(x))^2] = \mathbb{E}[(f(x) - f(x))^2] + \sigma^2 \end{aligned}$$

A partir de agora o erro é menor para aqueles que não eram
corridos de acordo com essas velhas ideias de probabilidade

$$Q_3 = np.multiply(Q_2, b_3) \quad \text{ou} \quad Q_2 = Q_3 * b_3$$

~ Depois temos que

$$\begin{aligned} \mathbb{E}[(f(x) - f(x))^2] &= \mathbb{E}[(f(x) - \mathbb{E}[f(x)] + \mathbb{E}[f(x)] - f(x))^2] \\ &= \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2] + (\mathbb{E}[f(x)] - \mathbb{E}[f(x)])^2 - 2(\mathbb{E}[f(x)] - \mathbb{E}[f(x)])(f(x) - \mathbb{E}[f(x)]) \end{aligned}$$

$$\begin{aligned} &= \mathbb{E}[(\mathbb{E}[f(x)] - f(x))^2] + (f(x) - \mathbb{E}[f(x)])^2 - 2\mathbb{E}[f(x)](\mathbb{E}[f(x)] - \mathbb{E}[f(x)]) \\ &\quad \underbrace{\mathbb{E}[f(x)] - f(x)}_{\text{bias}[f(x)]} + \underbrace{\mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2]}_{\text{var}[f(x)]} - 2(\mathbb{E}[f(x)] - \mathbb{E}[f(x)])(f(x) - \mathbb{E}[f(x)]) \\ &= \mathbb{E}[(f(x) - f(x))^2] = \text{bias}^2[f(x)] + \text{var}[f(x)] \end{aligned}$$

Since

$$\text{MSE} = \mathbb{E}[(y - f(x))^2] = \mathbb{E}[(f(x) - f(x))^2] + \sigma^2$$

and we've just shown that

$$\mathbb{E}[(f(x) - f(x))^2] = \text{bias}^2[f(x)] + \text{var}[f(x)]$$

hence, we have

$$\text{MSE} = \text{bias}^2[f(x)] + \text{var}[f(x)] + \sigma^2$$

Bias-Variance Equation

Let $E[MSE]$ be the expected mean-squared error of the fixed set of test instances over different samples of training data sets.

$$E[MSE] = \text{Bias}^2 + \text{Variance} + \text{Noise} \quad (2)$$

- In linear models, the bias component will contribute more to $E[MSE]$.
- In polynomial models, the variance component will contribute more to $E[MSE]$.

We have a trade-off, when it comes to choosing model complexity!

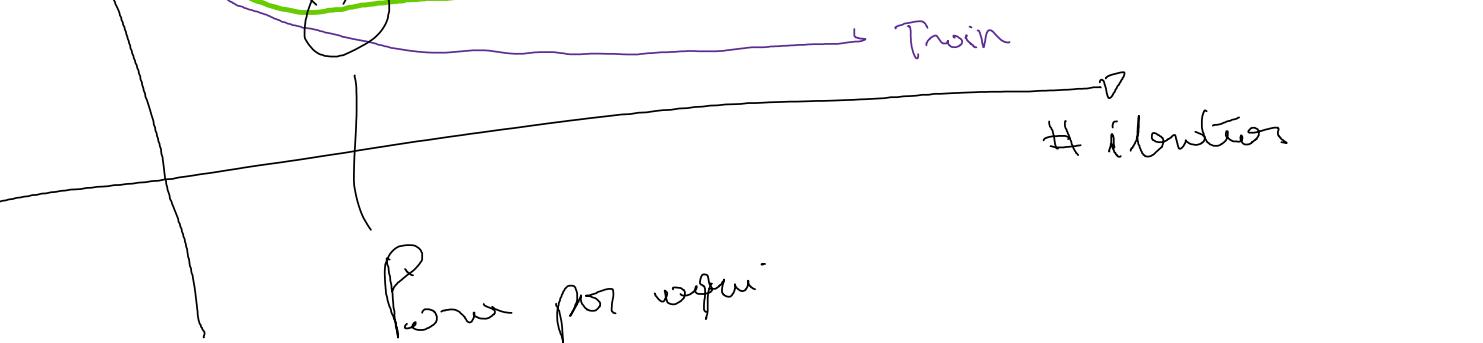
• Data Augmentation = Regra mais dura, quanto maior a quantidade

menor a chance de overfit

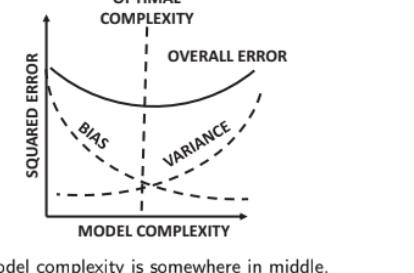
Você pode manipular dados já existentes para aprender mais dados. Ex: gerar uma foto, dar zoom em fotos

$$\mathbb{E}[(f(x) - f(x))^2] = \text{bias}^2[f(x)] + \text{var}[f(x)]$$

• Early Stopping = Você pega o trein set e o test set para o ponto de treinamento onde a diferença entre os erros ficar grande



The Bias-Variance Trade-Off



Bias-Variance Trade-off: Setup

- Imagine you are given the true distribution \mathcal{D} of training data (including labels).
- You have a principled way of sampling data sets \mathcal{D}' from the training distribution \mathcal{D} .
- Imagine you create an infinite number of training data sets (and trained models) by repeated sampling.
- You have a set T of unlabeled test instances.
- The test set T does not change over different training data sets.
- Computer prediction of each instance in T for each trained model.

Informal Definition of Bias

- Compute averaged prediction of each test instance x over different training models $f_i(x)$.
- Averaged prediction of test instance will be different from true (unknown) model $f(x)$.
- Difference between (averaged) $f_i(x)$ and $f(x)$ caused by erroneous assumptions/simplifications in modeling \Rightarrow Bias
- Example: If the true (Unknown) model $f(x)$ were an order-4 polynomial, and we used any polynomial of order 3 or greater to fit \mathcal{D} , then would be bias.

Informal Definition of Variance

- The value $f_i(x)$ will vary with \mathcal{D} for fixed x .
- The prediction of the same test instance will be different over different trained models.
- All those predictions cannot be simultaneously correct \Rightarrow Variation contributes to error
- Variance of $f_i(x)$ over different training data sets \Rightarrow Model Variance
- Lower order model will have high variance.

Bias-Variance Equation

Let $E[MSE]$ be the expected mean-squared error of the fixed set of test instances over different samples of training data sets.

$$E[MSE] = \text{Bias}^2 + \text{Variance} + \text{Noise} \quad (3)$$

- In linear models, the bias component will contribute more to $E[MSE]$.
- In polynomial models, the variance component will contribute more to $E[MSE]$.

We have a trade-off, when it comes to choosing model complexity!

• J₁ Penalty

Baixar hidden units, de modo que leve a restrição mais estrita

→ Ensemble Methods

• Técnicas para reduzir o erro de generalização ao se combinar diversos modelos. → Estratégia de ML chamada model averaging

• Ponto das promotoras que modelos diferentes não são certos ou nem mesmo errados

• Você pode tentar diferentes pontos de um training set de o seu dataset para grande o suficiente

• Com isso, a variação pode ser reduzida. Hipoteticamente, com dados infinitos, elas chegarão a zero

• Você pega os mesmos ou a média e usa

• no final, em tipo de operação vai gerar resultados melhores que nem só

• modelo em todo training set

• modelo em todo training set

• modelo com poucas unidades neurais

• modelo com poucas unidades neurais</p

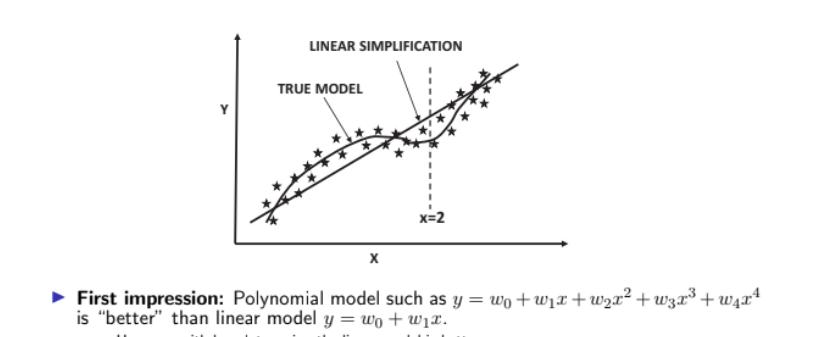
- The error on test data might be caused by several reasons.
 - Overfitting means might be low (overfitting), and poor convergence
 - Overfitting shows up as a large gap between in-sample and out-of-sample accuracy.
 - First solution is to collect more data.
 - New data might not always be available!

Penalty-Based Regularization

- Key techniques to improve generalization in NNs:
 - Penalty-based regularization
 - Constraints the shared parameters
 - Adding noise to input or output
 - Adding noise and stochasticity to input or hidden units

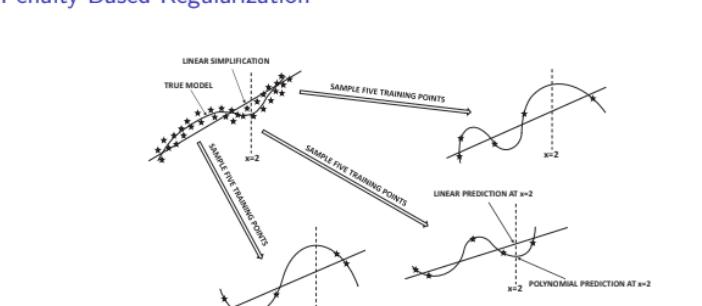
Penalty-Based Regularization

Revising Example: Predict y from x



- First important: Polynomial model such as $y = w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$ is "better" than linear model $y = w_0 + w_1x$.
- However, with less data, using the linear model is better.

Penalty-Based Regularization



- Zero error on training data but wildly varying predictions of $x = 2$

Economy in Parameters

- A lower-order model has economy in parameters.
 - Using a higher-order model has parameters, which might not use the parameters.
 - Economy in parameters decreases overfitting.
- Choosing a neural network with fewer layers per layer enforces economy.
- Reducing the number of parameters is a hard way.
- We can also penalize parameters in a soft way.

Summary

The Penn Machine Learning Toolkit
Penalty-Based Regularization
 L_2 Regularization and Soft Economy vs Hard Economy
 L_1 Regularization and Sparse Representations
Economy: Model Selection, Feature Selection, Sparse Representations
Regularization: Feature Selection and Averaging
Data Parallelism: Ensemble
Early Stopping
Cross Validation
Model Selection
What About Supervised Pre-training?

L_2 Regularization and Soft Economy vs Hard Economy

- Fixing the architecture up front is an inflexible solution.
- A softer solution uses a larger model but imposes a (tunable) penalty on parameter size.

$$g = \sum_{i=1}^n w_i x_i^2 \quad (4)$$
- Loss function: $L = \sum_{(x,y) \in D} (y - \hat{y})^2 + \lambda \sum_{i=1}^n w_i^2$
- Regularization
- The (tunable) value of λ decides the level of regularization.
- Softer approach with a complex model performs better!

Effect on Updates

- For learning rate α , effect on update is to multiply parameter with $(1 - \alpha\lambda) \in (0, 1)$.

$$w_i \leftarrow w_i(1 - \alpha\lambda) - \alpha \frac{\partial L}{\partial w_i}$$
- Interpretation: Decay-based forgetting
- Unless a parameter is important, it will have small absolute value.
- More weight decay means more than weight decaying or zero.
- A forcing gradient in the loss function from monitoring the training data, because only significant and repeated updates will be reflected in the weights.

Summary

The Penn Machine Learning Toolkit
Penalty-Based Regularization
 L_2 Regularization and Soft Economy vs Hard Economy
 L_1 Regularization and Sparse Representations
Economy: Model Selection, Feature Selection, Sparse Representations
Regularization: Feature Selection and Averaging
Data Parallelism: Ensemble
Early Stopping
Cross Validation
Model Selection
What About Supervised Pre-training?

L_1 -Regularization

- In L_1 -regularization, an L_1 -penalty is imposed on the loss function.

$$L = \sum_{(x,y) \in D} (y - \hat{y})^2 + \lambda \sum_{i=1}^n |w_i|$$
- Update has slightly different update equation at least for the case when $w_i \neq 0$.

$$w_i \leftarrow w_i - \alpha \lambda \text{sgn}(w_i) - \alpha \frac{\partial L}{\partial w_i}$$
- The value of a_i is the partial derivative of $|w_i|$ w.r.t. w_i .

$$a_i = \begin{cases} -1 & w_i < 0 \\ +1 & w_i > 0 \end{cases}$$

Weight space loss term $\mathcal{L}(y, \hat{y})$ and L_1 - and L_2 -regularization terms

L_1 - or L_2 -Regularization?

- L_1 -regularization leads to sparse parameter learning.
 - Zero values of w_i can be interpreted as inactive neurons in the neural network.
- L_2 -regularization generally provides better performance.
- L_2 is differentiable and can be used in different techniques.

Summary

The Penn Machine Learning Toolkit
Penalty-Based Regularization
 L_2 Regularization and Soft Economy vs Hard Economy
 L_1 Regularization and Sparse Representations
Economy: Model Selection, Feature Selection, Sparse Representations
Regularization: Feature Selection and Averaging
Data Parallelism: Ensemble
Early Stopping
Cross Validation
Model Selection
What About Supervised Pre-training?

Bagging and Sub-sampling
Bagging and Sub-sampling are imperfect simulations of drawing the training data from a base distribution.
Nevertheless, the variance of this approach is still lower than that of constructing a single model on the entire training data set.

Connections with Noise Injection

- L_2 -regularization with parameter λ is equivalent to adding Gaussian noise with variance λ^{-1} to input.
- Regularization terms of noise will be minimized with simple models (under parameters).
- Result is only true for single layer networks (linear regression).
- Mean value of noise is providing generalization.
- Similar results can be shown for denoising autoencoders.

Summary

The Bias-Variance Trade-off
Generalization Error in Model Training and Evaluation
Penalty-Based Regularization
Lasso Regression and Elastic Net, Soft-Elastic Net vs Hard-Elastic Net
Penalizing Hidden Units: Learning Sparse Representations
Feature Selection
Parameter Selection
Penalized Least Squares
Penalized Maximum Likelihood
Penalized Model Selection and Averaging
Cross-Validation
Data Partitioning
Cross-Validation
Data Partitioning
Out-of-Sample Performance
What About Sparse Representations?

Penalizing Hidden Units

- One can also penalize hidden units (autoencoders).
- Applying L_2 -penalty leads to sparse activations.
- Depth-forward modification of backpropagation.

Penalty contributions from hidden units are added up in backward phase.

$$E' = J + \lambda \sum_i |h_i| \quad (5)$$

Where:
► J is the total number of units in the network.
► h_i is the value of the i th hidden unit.
► λ is the regularization parameter.

Ensemble Methods

- Inspired in Bias-Variance trade-off.
- Try to reduce either the bias or the variance without affecting the other component.
- Ensemble methods are commonly used in Machine Learning.
- Two examples:
 - Bagging - Variance reduction.
 - Boosting - Bias reduction.

Ensemble Methods

- Most ensemble methods in NN are focused on variance reduction.
- This is because neural networks are valued for their ability to build arbitrarily complex models with relatively low bias.
- But arbitrarily complex models lead to **OVERTFITTING**.

Therefore, the goal of most ensemble methods in NN is variance reduction.

Summary

The Bias-Variance Trade-off
Generalization Error in Model Training and Evaluation
Penalty-Based Regularization
Lasso Regression and Elastic Net, Soft-Elastic Net vs Hard-Elastic Net
Penalizing Hidden Units: Learning Sparse Representations
Feature Selection
Parameter Selection
Penalized Least Squares
Penalized Maximum Likelihood
Penalized Model Selection and Averaging
Cross-Validation
Data Partitioning
Cross-Validation
Data Partitioning
Out-of-Sample Performance
What About Sparse Representations?

Bagging and Sub-sampling

- If a sufficient number of samples is available, after all, the variance of most types of output predictions can be asymptotically reduced to zero.
- One approach: Predictions are repeated and averaged over training data sets.
- With a sufficient large number of training data sets is used, the variance can be reduced to zero (infinite source of data).
- Although the variance of individual predictions is high, the variance of the aforementioned methodology has still better variance characteristics than a single execution of the model on the entire training data set.

Bagging and Sub-sampling

- The predictions on a particular test instance, which are obtained from the models built with different training sets, are then averaged to create the final prediction.
- One can average either the real-valued prediction (e.g., probability estimates of class membership) or the predictions.
- In the case of real-valued predictions, better results are sometimes obtained by using the median of the values.

Bagging and Sub-sampling

- Bagging:
 - Sample with replacement.
 - Sample size x (commonly $x = m$, in which case we have sample size equals training set size).
 - When $x = m$, we'll have duplicates and about a fraction $(1 - 1/m)^m \approx 1/e$ of the original data set will not be included.
 - Repetition of the same x leads the k models to a given test instance.
 - Average the results to yield a single robust prediction.
 - Common: $x = m$; best results with $x < m$.

Bagging and Sub-sampling

- Sub-sampling:
 - Sample without replacement.
 - Predictions are averaged.
 - Essential to choose $x < m$, since $x = m$ would imply same training data set and identical results.

Bagging and Sub-sampling

- All the variance cannot be removed by using bagging or sub-sampling.
- Predictions of test instances from different samples will be positively correlated.
- The average of a set of random variables that are positively correlated will always have a variance that is proportional to the level of correlation.

Bagging and Sub-sampling
Bagging and Sub-sampling are imperfect simulations of drawing the training data from a base distribution.
Nevertheless, the variance of this approach is still lower than that of constructing a single model on the entire training data set.

- Main challenge is to construct multiple training models → highly inefficient.
- But it can be fully parallelized → use with GPUs.

Summary

The Five-Fold Cross-Validation
Generalization Issues in Model Tuning and Evaluation
Ensemble Methods
• Bagging
• Random Forests
• AdaBoost
• Stochastic Gradient Descent
• Gradient Boosting
• Feature Selection
• Feature Extraction
• Feature Learning
• Feature Representations
• Ensemble Methods
• Parametric Model Selection and Averaging
• Data Parallelism
• Data Parallelism Ensemble
• Deep Learning
• Deep Learning Pre-training
• Deep Model Selection and Pre-training

Parametric Model Selection and Averaging

- Problems:
 - Large number of hyper-parameters and configurations.
 - Hard to find good configuration.
- Strategy: Hold out a portion of the training data.
- Select the model out of the pool providing highest performance.

Summary

The Five-Fold Cross-Validation
Generalization Issues in Model Tuning and Evaluation
Ensemble Methods
• Bagging
• Random Forests
• AdaBoost
• Stochastic Gradient Descent
• Gradient Boosting
• Feature Selection
• Feature Extraction
• Feature Learning
• Feature Representations
• Ensemble Methods
• Parametric Model Selection and Averaging
• Dropout
• Data Parallelism Ensemble
• Early Stopping
• Deep Learning
• Deep Learning Pre-training
• Deep Model Selection and Pre-training

Basic Dropout Training Procedure

- For each training instance do:
 - Sample each node in the network in each layer (except output layer) with probability p .
 - Keep only edges for which both ends are included in network.
 - Retain all edges for which both ends are not included in network.
- Note that weights are shared between different sampled networks.
- A different neural network is used for every small mini-batch of training examples.
- The number of NNs is rather large in Dropout.

Basic Dropout Testing Procedures

- First procedure
 - Perform repeated sampling (like monte carlo) and average results.
 - Generally sampling will be problematic; subject averaging to bootstrap.
- Second procedure with weight scaling inference rule (more common)
 - Assume that each node has probability p of being sampled.
 - Perform single inference or full network with down-scaled weights.

Why Does Dropout Help?

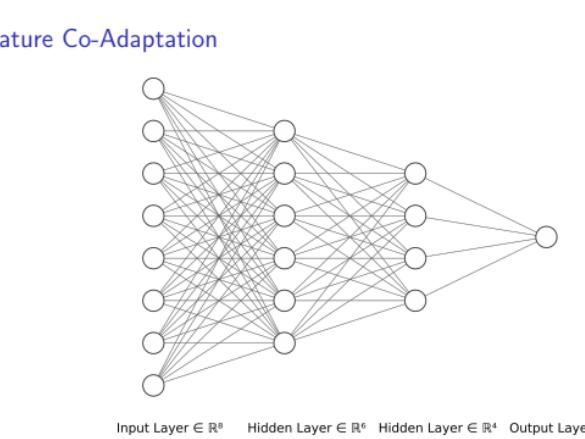
- By dropping nodes, we are forcing the network to learn without the presence of some inputs (in each layer).
- Will resist co-adaptation, unless the features are truly synergistic.
- Will create many (smaller) groups of self-sufficient predictors.
- Many groups of self-sufficient predictors will have a model-averaging effect.

The Regularization Perspective

- One can view the dropping of a node as the same process as adding masking noise.
 - Noise is added to both input and hidden layers.
- Adds to the loss function via regularization.
- Forces the weights to become more spread out.
 - Weights are distributed across weights based on sampling.

Practical Aspects of Dropout

- Typical dropout rate (i.e., probability of exclusion) is somewhere between 20% to 50%.
- Better to use a larger network with Dropout to enable learning of independent representations.
 - Dropout is applied to both input and hidden layers.
 - Large learning rate variance and large momentum.
 - Impose a max-norm constraint on the size of network weights.
 - Norm of input weights to a node upper bounded by constant.



One-Way Adaptation

- Consider a single-hidden layer neural network.
 - All input units and one of the hidden nodes are fixed to random values.
 - But the features will adapt to the other half (random features).
- Feature co-adaptation is useful when rate of training varies across different parts of network over time.
 - Facilitate a mechanism of training efficiency (over and above the spring).

Why is Feature Co-Adaptation Bad?

- We want features working together only when essential for prediction.
 - If we have features adapting to each other because of inefficiencies in training, then we are not using features adequately to make better predictions.
 - Due to co-adaptation, the features are not able to learn effectively.
- We want many groups of minimally essential features for robust prediction => Feature co-adaptation.
 - We do not want a few large and inefficiently created groups of co-adapted features.

Feature Co-Adaptation

- The process of training a neural network often leads to a high level of dependence among features
- Different parts of the network train at different rates
 - Some parts of the network tend to adapt faster than others.
- This is referred to as feature co-adaptation.
- Models trained this way are sensitive to nuances of specific training data => **OVERFITTING**

Summary

The Five Detector Train OFF
Generalization Issues in Model Tuning and Evaluation
Early Stopping
Data Augmentation
Ensemble Methods
Data Perturbation Ensembles
Data Shifting
Unsupervised Pre-training
What Makes Supervised Pre-training?

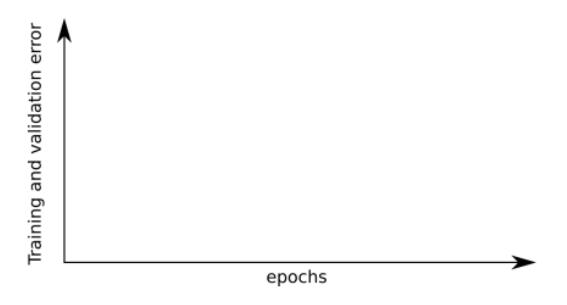
Data Perturbation Ensembles

- Most of the ensemble techniques discussed so far are either sampling-based ensembles or model-centric ensembles
- Dropout can be considered an ensemble that adds noise to the data in an indirect way
- Simplest case = noise is added to input data and weights are trained on the distorted data
- Repeat the process and average results.
- This is the generalizable method, not specific to neural networks.
- It is simple and noisy, however, it must be carefully calibrated.
- Dropout indirectly adds noise to hidden layers by randomly dropping nodes

Data Perturbation Ensembles

- Data augmentation can often greatly improve the accuracy of a learner by adding new training examples.
- But they are not perturbation schemes because the augmented examples are created with a calibrated procedure and understanding of the domain at hand.

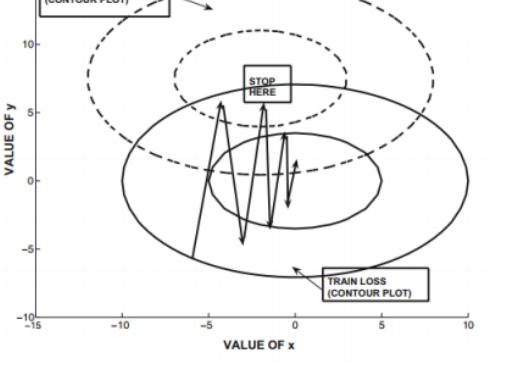
Early Stopping



Early Stopping

- Motivation:
 - We use optimization to train NNs until convergence.
 - Optimizes the loss on training data.
 - Not necessarily on the out-of-sample test data.
 - Final steps cause overfitting and generalization problems.
 - Almost always used.

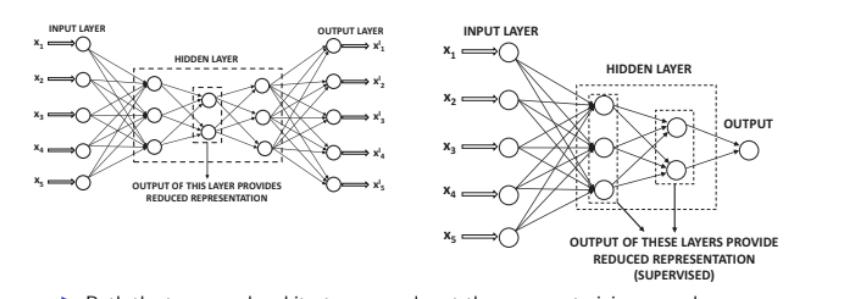
Early Stopping



Unsupervised Pre-training – Importance of Initialization

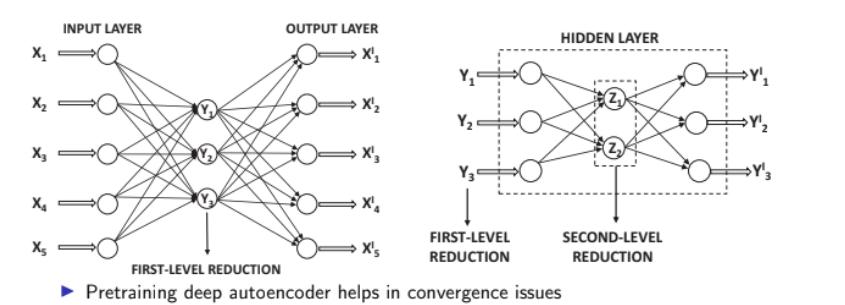
- Bad initializations can lead to unstable convergence.
- Typical approach is to initialize to a Gaussian with variance $1/r$, where r is the degree of freedom of the layer.
- Xavier initialization uses both in-degree and out-degree.
- Pretraining goes beyond these simple initializations by using the training data.

Types of Base Applications



- Both the two neural architectures use almost the same pretraining procedure

Layer-Wise Pretraining a Deep Autoencoder



Summary

The Five Detector Train OFF
Generalization Issues in Model Tuning and Evaluation
Early Stopping
Data Augmentation
Ensemble Methods
Data Perturbation Ensembles
Data Shifting
Unsupervised Pre-training
What Makes Supervised Pre-training?

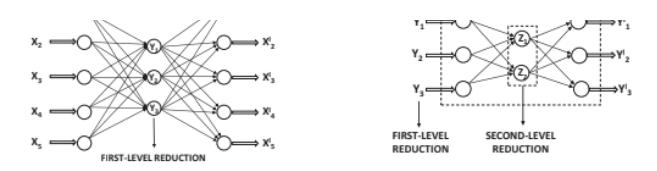
Types of Pretraining

- Unsupervised pretraining: Use training data **without** labels for initialization. Improves convergence behavior.
- Supervised pretraining: Use training data **with** labels for initialization. Labels are compared via loss terms.
- Focus on unsupervised pretraining.

Pretraining a Supervised Learner

- For a supervised learner with k hidden layers:
 - Remove output layer and create an autoencoder with $(k-1)$ hidden layers.
 - Pretrain first layer with unlabeled data.
 - Keep only weights from encoder portion and cap with output layer.
 - Pretrain remaining layers.
 - Fine-tune all layers.

INPUT LAYER HIDDEN LAYER OUTPUT LAYER



Why Does Pretraining Work?

- ▶ Pretraining (step 1) brings the activations of the neural network to the manifold of the data distribution.
- ▶ Features correspond to repeated patterns in the data.
- ▶ Fine-tuning learns to combine/modify relevant ones for inference.
 - Reuse: -72 neurons share first layer in generic images

Thank you!
tveret@ic.uff.br