

Bias, variance and regularization

Jan Chorowski
Instytut Informatyki
Wydział Matematyki i Informatyki Uniwersytet
Wrocławski
2020

Where are we?

Goal:

do well

- low error rate
- proper fraction
- essentially best \mathcal{L}

on unseen data

- low error rate on test data
- proper regularization
- the difference between \mathcal{ML} and OPT.
- goal for today

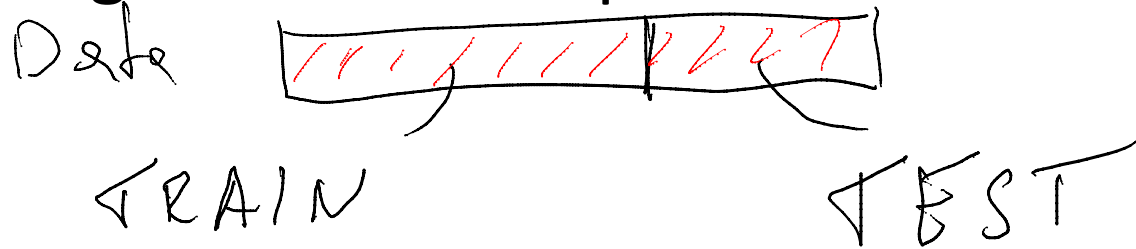
Generalization, the goal of learning

- Problem:
 - We care about the performance on **all the data**
 - We have only a **training sample**
 - **Under-fitting:**
Doesn't work even on **training samples**
Too weak classifier (cannot express the relation in the data), bad training ...
 - **Over-fitting:**
Too powerful classifier will perfectly interpolate the training data (even the noise in it!) and do poorly on **unseen samples**
- Questions:
 - How to estimate the **generalization** (performance on all data)? -> Honest estimates
 - How to control the capacity of a model?
 - Can we provably ensure good generalization performance -> Learning Theory

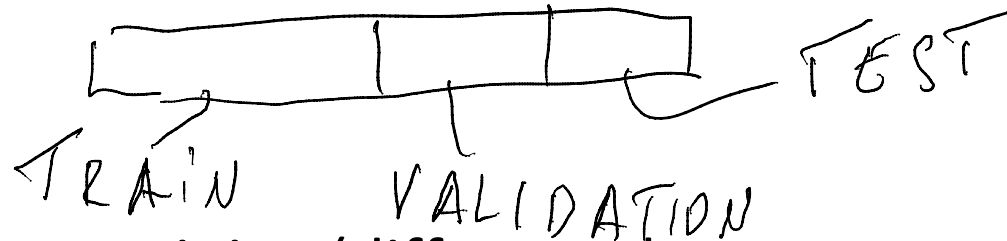
Honest estimates: Hold-out set

Large data case!!!

- Split the training data into two parts:



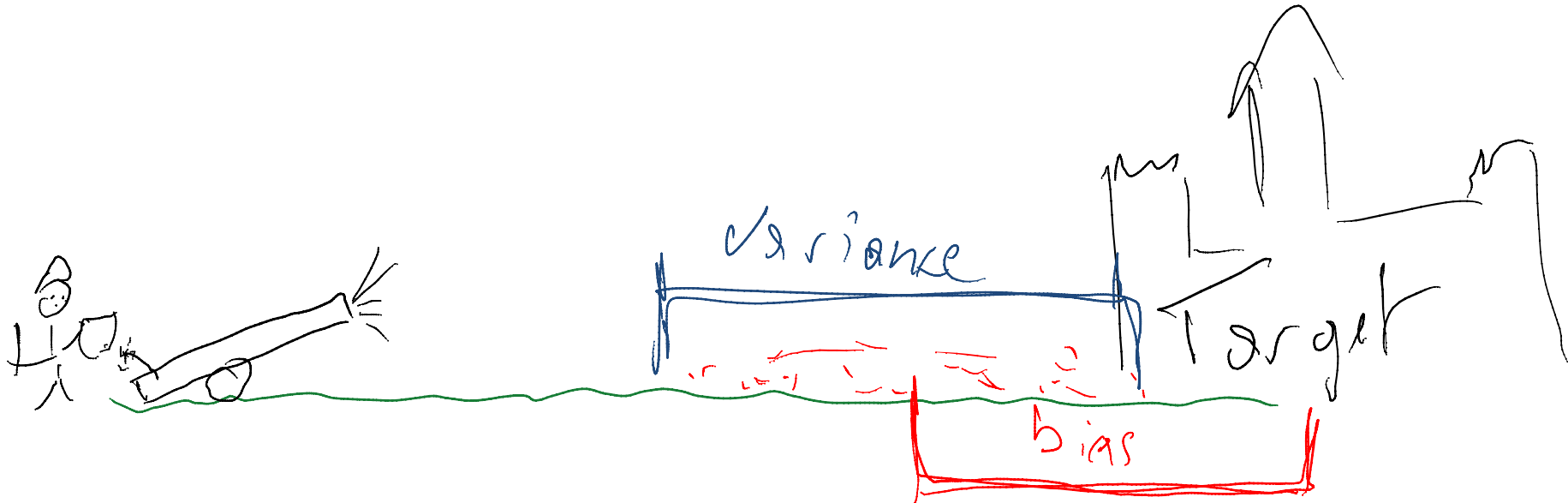
- Train only on training, then test on testing.
- Often we do a three-way split:



- Then:
 - Train many models on training (different algos, parameters)
 - Use validation to choose best model
 - Test on testing
- Other techniques: Cross-Validation and bootstrap

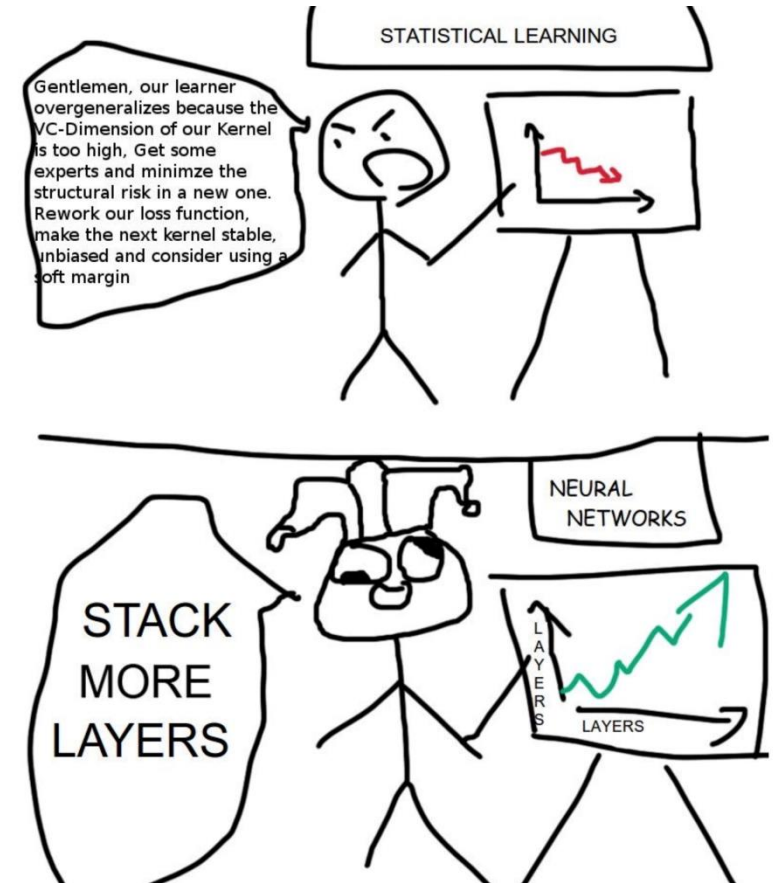
Bias-Variance: two sources of error!

- The **bias** captures how well our family of functions (hypothesis space) matches the data.
Intuition: systematic error
Large neural nets have a low bias (universal approximation!)
- The **variance** captures how results of different training runs vary
Deep Learning is unstable: random init, SGD...



How to lower the bias?

- Choose more powerful/better models:
 - **More hidden neurons**
 - Use better network architecture
 - Understand the data and choose a matching model
 - Describe the data with more attributes
 - Better data transformation
- This usually increases the Hypothesis space



How to lower variance?

- **Get more data**
... or generate synthetic, e.g. rotate and shear pictures
- **Constrain the models:**
 - Simpler models
 - Shorter training (early stopping)
 - Regularize the models
 - Select only the most important inputs
- **Average the models (very powerful)**
 - Also called “ensemble learning”, boosting, bagging
 - Requires that the models make uncorrelated errors
 - Some neural net techniques have an ensembling interpretation:
 - Dropout: ensemble of many nets that share weights, but use different neurons
 - ResNets: sum of many shallow nets formed by subsets of layers

Classical Regularization

- Limit the size of parameters
- Ridge regression: $\min_{\Theta} \sum_i (\Theta^T x^{(i)} - y^{(i)}) + \frac{\lambda}{2} \sum_j \Theta_j^2$
- Interpretation:
 - Unique solution
 - Small $\frac{\partial \Theta^T x}{\partial x}$
 - Use all features to decide (L2 reg), select features (L1 reg)

Weight Decay == Ridge Regression in DL

The intuitions:

- Start with many weights (larger nets train easier!)
- Don't depend on a single net. How?
- Force all the weights to decrease
- Subtract a little bit in each training iteration:
$$\Theta \leftarrow \Theta - \alpha(\nabla_{\Theta} (J) + \beta\Theta) = \Theta - \alpha\nabla_{\Theta} J - \alpha\beta\Theta \quad (\text{weight decay})$$
- Note: this minimizes $J(\Theta) + \frac{\beta}{2} \sum_j (\Theta_j)^2$
- Note: usually you don't decay the biases

Probabilistic interpretation of Weight Decay

- The gradient step:

$$\Theta := \Theta - \alpha(\nabla_{\Theta} (J) - \beta\Theta)$$

Corresponds to minimizing:

$$J(\Theta) + \frac{\beta}{2} \sum_j (\Theta_j)^2$$

Now try to find a probabilistic interpretation!

Bayesian approach

1. Make some models more probable than others
2. Set a **prior** probability distribution over Θ
3. For example:
 - weights are normally distributed
 $p(\Theta_i) \sim \mathcal{N}(0, \sigma_{\Theta})$
4. Treat Θ as a random variable too
 $P(Y|X, \Theta)$ i.e. y depends on x and Θ which is randomly sampled too

Bayes theorem

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$

Interpretation: how our estimate of A changes after seeing B .

Why?

$$p(A, B) = p(A|B)p(B) = p(B|A)p(A)$$

Then divide by $p(B)$

Bayesian approach to ML

- What is the model probability after seeing the examples in Data \mathcal{D} ?

$$p(\Theta|\mathcal{D}) = \frac{p(\mathcal{D}|\Theta)p(\Theta)}{p(\mathcal{D})}$$

How to make predictions? Integrate over all models:

$$p(y|x, \mathcal{D}) = \int_{\Theta} p(y|x, \Theta)p(\Theta|\mathcal{D})d\Theta$$

Then

$$E[y|x, \mathcal{D}] = \int_y yp(y|x, \mathcal{D})dy$$

But computing $p(y|x, S)$ is often intractable :(

Maximum-a-posteriori

- Instead of predicting integrating over all Θ
- Use the maximally probable Θ :

$$\begin{aligned}\Theta_{MAP} &= \arg \max_{\Theta} p(\Theta | \mathcal{D}) \\ &= \arg \max_{\Theta} \left(\prod_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} p(y^{(i)} | x^{(i)}, \Theta) \right) p(\Theta)\end{aligned}$$

- It's like Max. Likelihood with the extra term.

Gaussian model MAP

$$\arg \max_{\Theta} \prod_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} p(y^{(i)} | x^{(i)}, \Theta) p(\Theta) =$$
$$\arg \max_{\Theta} \sum_{(x^{(i)}, y^{(i)}) \in \mathcal{D}} \log p(y^{(i)} | x^{(i)}, \Theta) + \log(p(\Theta))$$

Now if Θ_j are Gaussian with zero-mean:

$$p(\Theta_j) = \frac{1}{\sigma_{\Theta} \sqrt{2\pi}} e^{-\frac{\Theta_j^2}{2\sigma_{\Theta}^2}}$$

Then we recover the weight decay term:

$$-\log p(\Theta_j) \propto \Theta_j^2$$

Weight decay and momentum training

Recall that in weight decay:

$$TotLoss = Loss + \lambda \sum_i \Theta_i^2$$
$$\frac{\partial TotLoss}{\partial \Theta} = \frac{\partial Loss}{\partial \Theta} + \frac{\lambda}{2} \sum_i \Theta_i$$

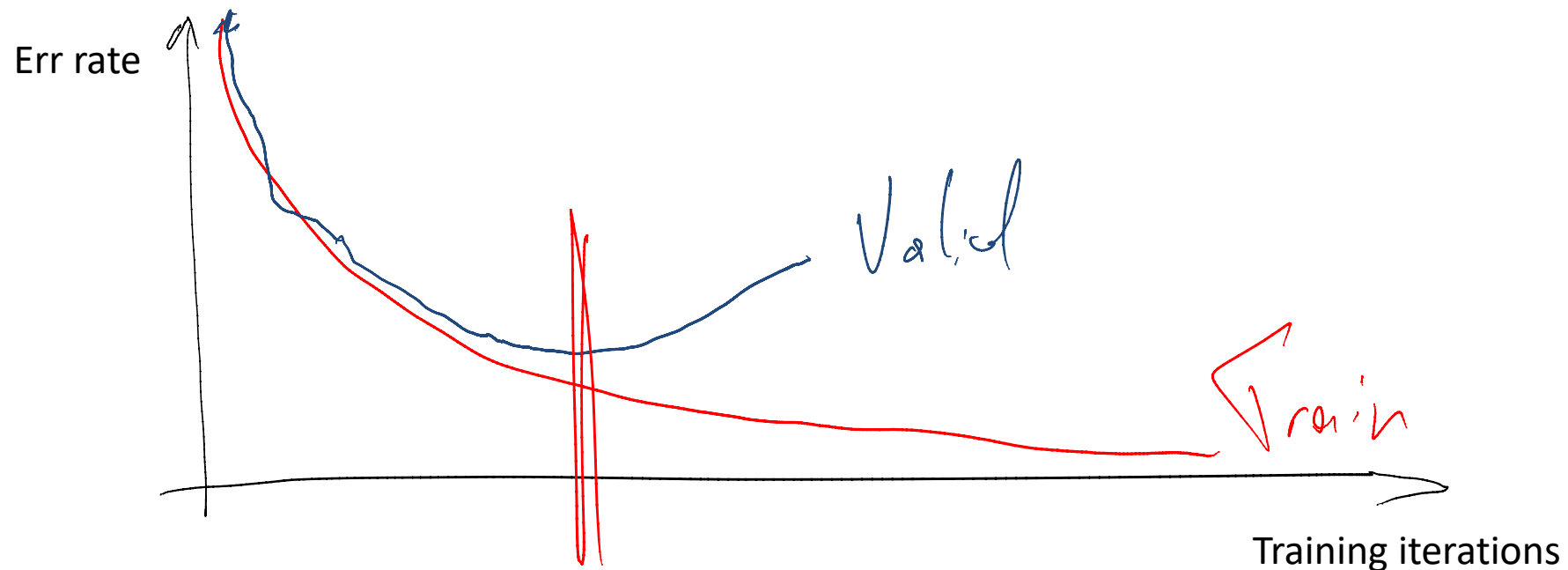
The **weight decay** term is stable over time -> gets boosted by momentum

Hacky trick (promoted by fast.ai): don't include weight decay in loss and gradient computation.

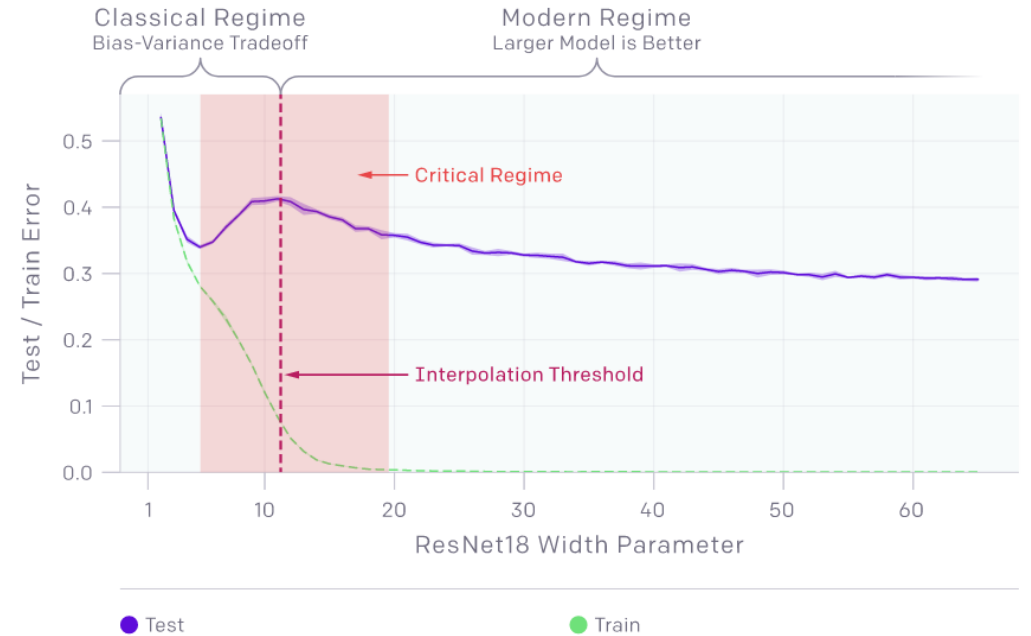
However, apply the weight decay step separately after the main optimizer.

Early stopping

- The net starts with small weights (we initialize it like that)
- As training progresses the weights grow (net specializes)
- At some point, it over-specializes
- Look for that moment, by monitoring a validation error!



With early stopping, large nets don't overfit



Weight norm constrains

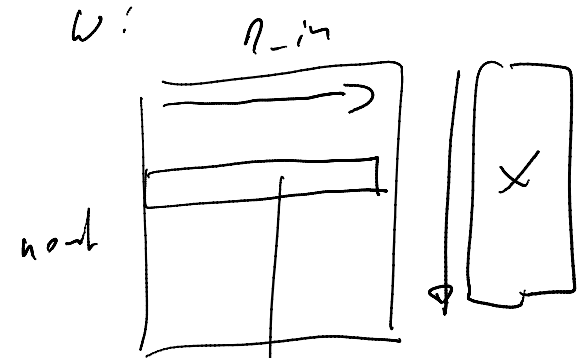
In Weight Decay we want small weights ...and pull all weights to 0!

Similar idea: constrain weights to be small, i.e.

$$\min_{W:} Wx$$

$$\forall i \text{ norm}(W_{i:}) < \delta$$

↳ constrain all rows.



↳ each row: weights of 1 neuron

How?

- 1 make an SGD update
- 2 constrain

$$W := \max\left(1, \frac{\text{sqrt}(\text{sum}(W \times W, 1))}{\delta}\right)$$

Weight noise

Weights can be large, but have limited precision
(thus it doesn't matter if they change a bit).

How? Add noise in SGD:

$$\Theta \leftarrow \frac{\partial \text{Loss}(\Theta + \text{noise})}{\partial \Theta}$$

for each minibatch:

$\Theta' = \Theta + \text{noise}$

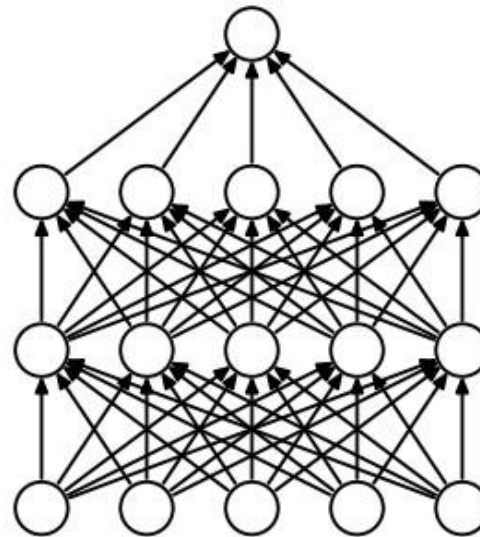
do fprop using Θ'

do bprop

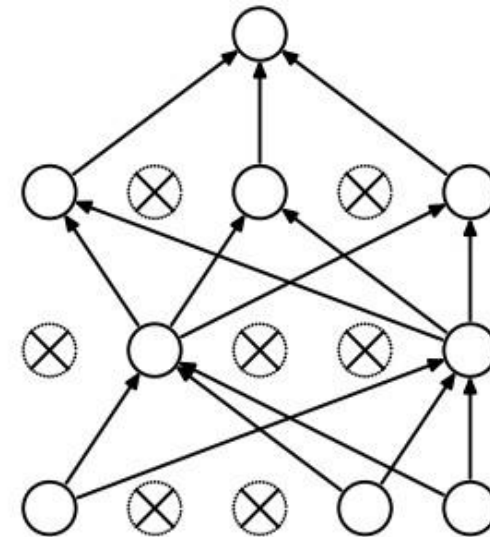
apply gradient.

Dropout

- At train time:
 - For each example, select with probability p which neurons will be used (not dropped out).
 - Multiply the outputs of other neurons by $1/p$ to compensate
- At test time:
 - Use all neurons
- Interpretation:
 - Prevents co-adaptation of neurons, it is harder for neurons to cooperate if any can be dropped-out
 - Trains infinitely many networks, each sharing selected neurons with the other ones



(a) Standard Neural Net



(b) After applying dropout.

Networks are overconfident

- We usually train with cross-entropy loss

$$L = \sum_{c=1}^{C_i} -p(y^{(i)} = c) \log p(y = c | x^{(i)})$$

- In the train data the class is certain, loss simplifies to

$$L = -\log p(y = y^{(i)} | x^{(i)})$$

- When model is 99% accurate...
- The only way to reduce loss is to make $p(y = y^{(i)} | x^{(i)}) = 1$
- This promotes overconfidence $\underline{p(\text{first guess})} \gg \underline{p(\text{second guess})}$

Label Smoothing

- Introduced in Inception V2 (arXiv:1512.00567)
- Assume train labels are more ambiguous:

$$L = \sum_c -p(y^{(i)} = c) \log p(y = c | x^{(i)})$$

- $p(y^{(i)} = c)$ is a smoothing distribution, e.g.

$$p(y^{(i)} = c) = \begin{cases} \beta, & \text{on correct class} \\ \frac{1 - \beta}{C - 1}, & \text{otherwise} \end{cases}$$

- Even better: smooth the $1 - \beta$ according to class marginal probabilities

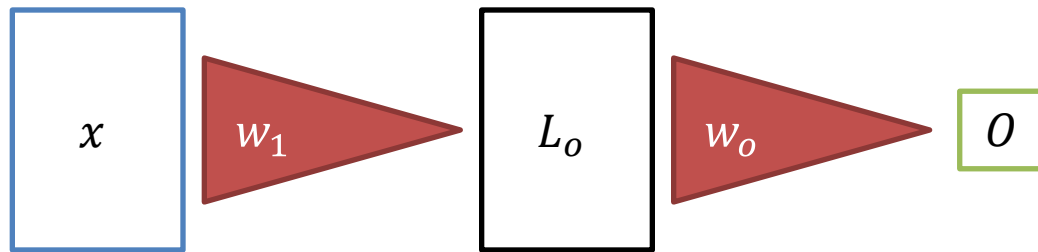
Effects of Label Smoothing

- Reduces overconfidence and regularizes
- Also prevents gradient vanishing:
 - Without smoothing SoftMax derivative is $p_{\Theta}(Y_i|X_i) - [Y_i = c]$
 - This vanishes when $p_{\Theta}(Y_i|X_i) \approx 1$
 - Effectively the model stops training on correctly classified instances

Label Smoothing vs Other Regularizers

At a high level, all regularizers want to forbid large changes of output for small changes of input.

- E.g. weight decay



Magnitude of w_o controls the output sensitivity $\frac{\partial o}{\partial L_o} = w_o^T$

- Label smoothing may be easier to use:
 - Easy to say how smooth the output should be
 - Hard to say how large the weights should be

Mixup

Mixup (<https://arxiv.org/pdf/1710.09412.pdf>)

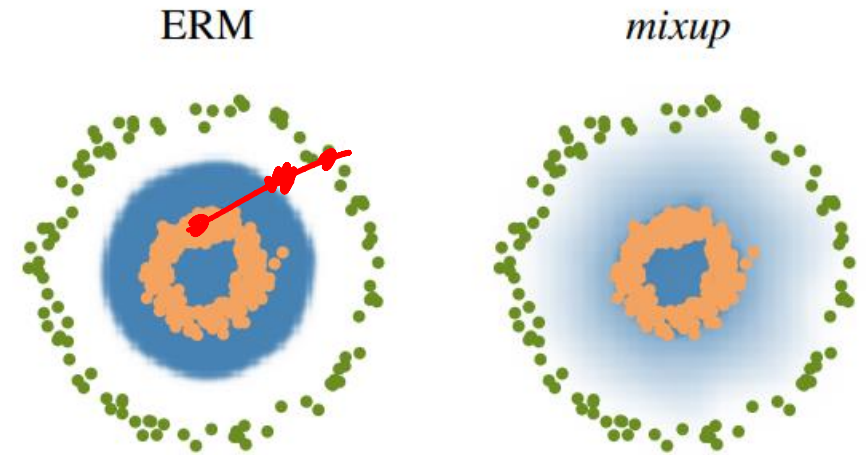
- Take two data samples:
 $(x^{(i)}, y^{(i)})$ and $(x^{(j)}, y^{(j)})$.

- Make a *new* sample

$$\tilde{x} = \lambda x^{(i)} + (1 - \lambda)x^{(j)}$$

$$\tilde{y} = \lambda y^{(i)} + (1 - \lambda)y^{(j)}$$

- Intuition: again, this forces the network's output to gradually change



Summary of NNet regularizations

- Choose proper architecture
add or remove neurons or whole layers
 - Share weights between neurons
 - Example: convolutional networks
- Use weight decay or other weight priors
L2 (Gaussian), or L1 (Laplace)
- Use Early stopping
Monitor validation error as training progresses. Stop when it starts to increase
- Use dropout -> randomly remove some neurons
- Use weight noise

Hyperparameters: Art or Science

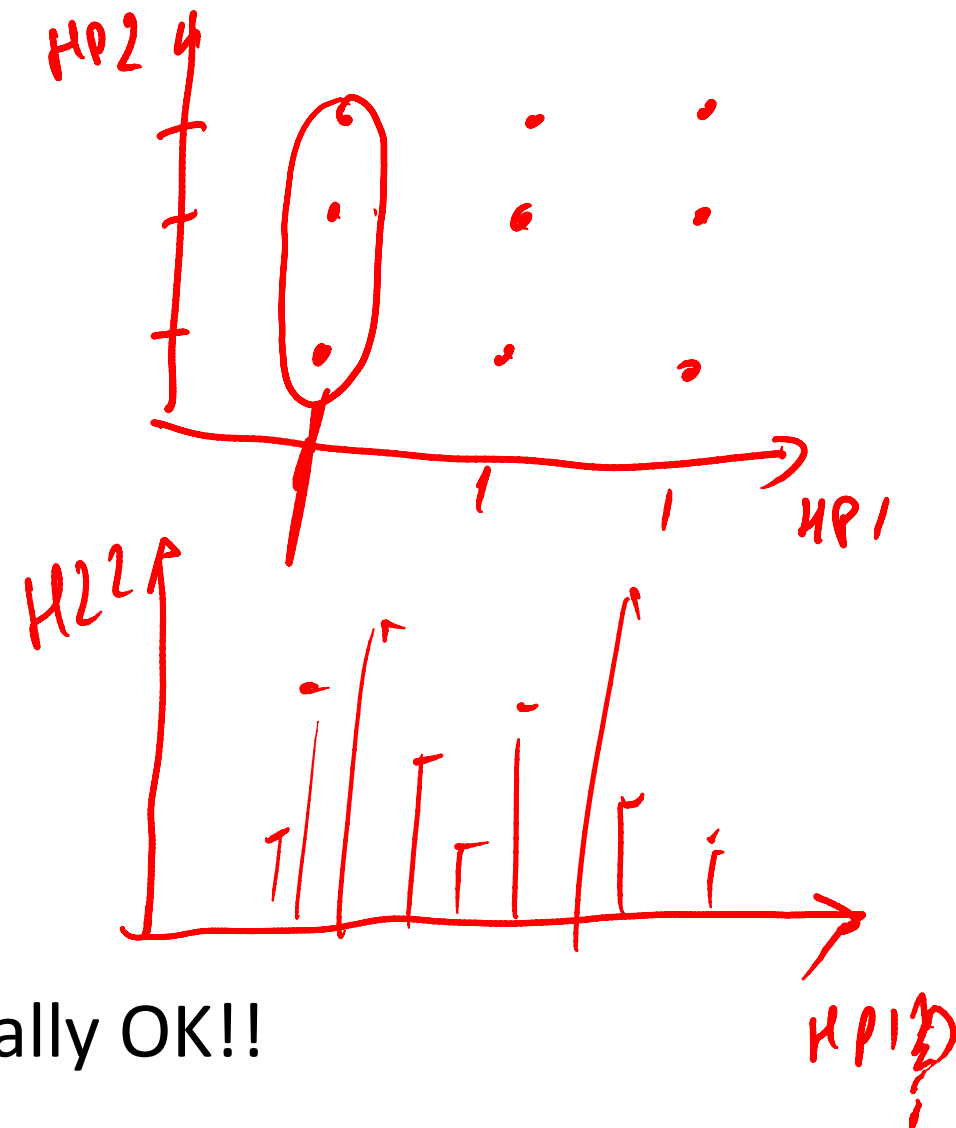
Training a net requires many decisions:

- How many layers, how many neurons
- What weight decay
- What learning rate, what momentum
- What dropout
- Use Batch norm?

We choose them on the validation set!

Q: how to organize search?

Random hyperparam search is easy and usually OK!!



Hyperparameters: smarter approaches

The random search is surprisingly good.

However, more sophisticated approaches exist.

- <https://github.com/JasperSnoek/spearmint> provides a method that uses Gaussian Processes to estimate hyperparameter impact
- Many cloud providers have their own tools, e.g. <https://cloud.google.com/ai-platform/training/docs/using-hyperparameter-tuning>