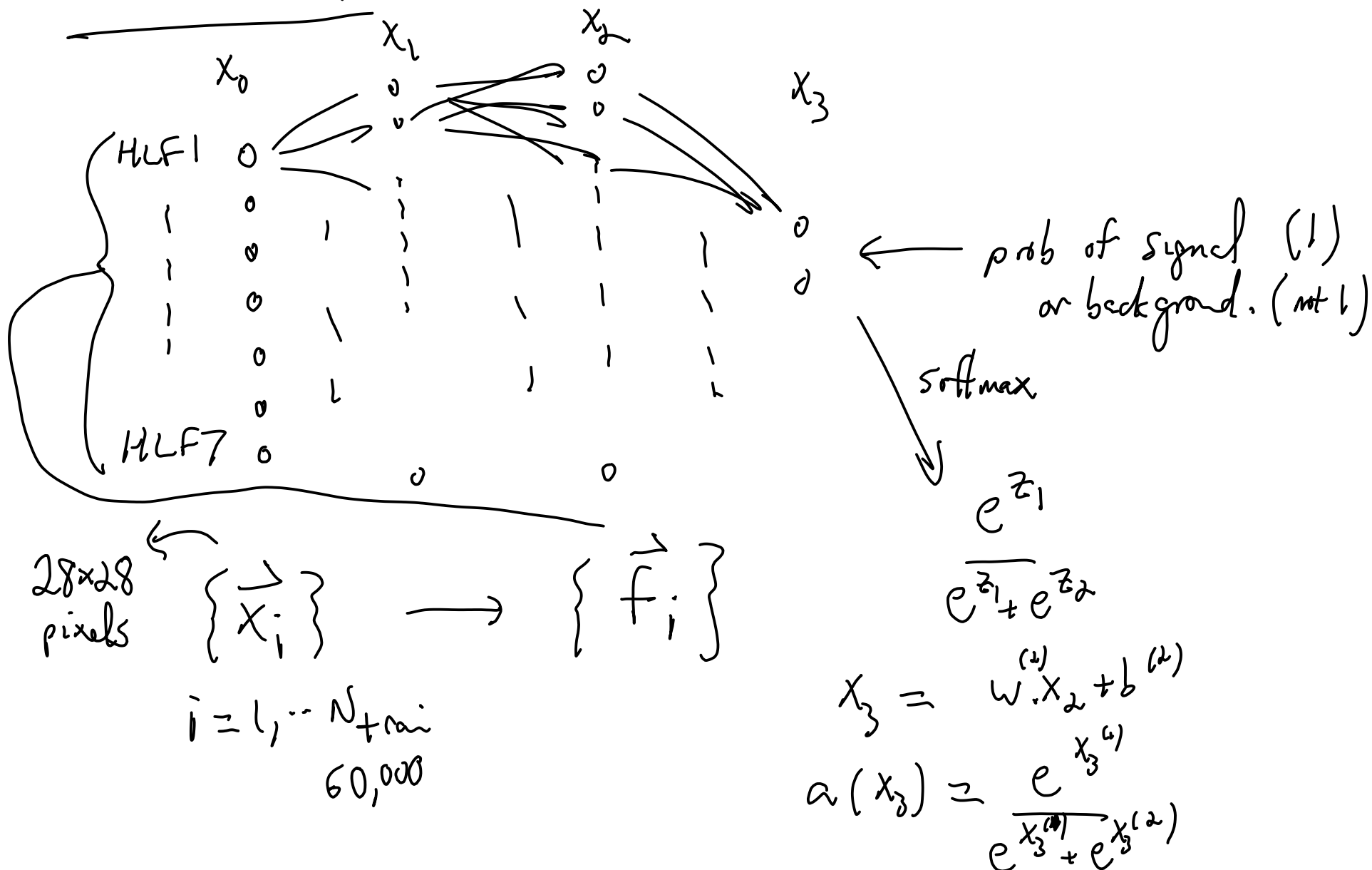


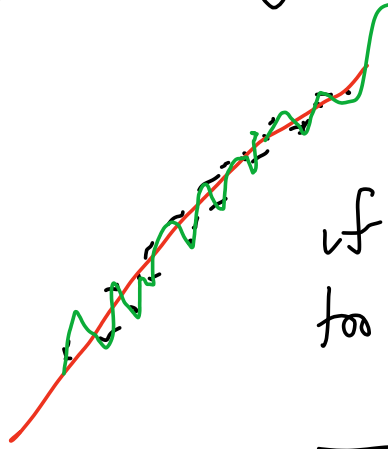
Friday, January 29, 2021 3:05 PM

Lecture 4



Friday, January 29, 2021 3:10 PM

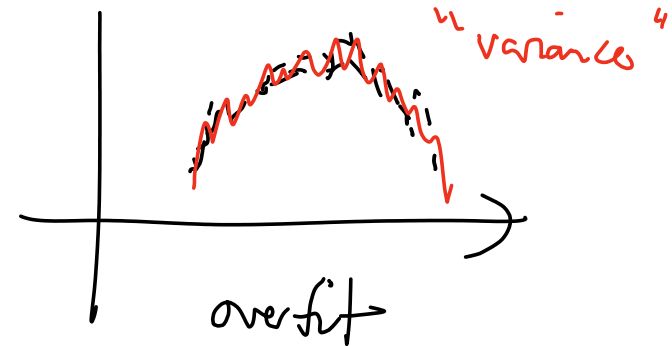
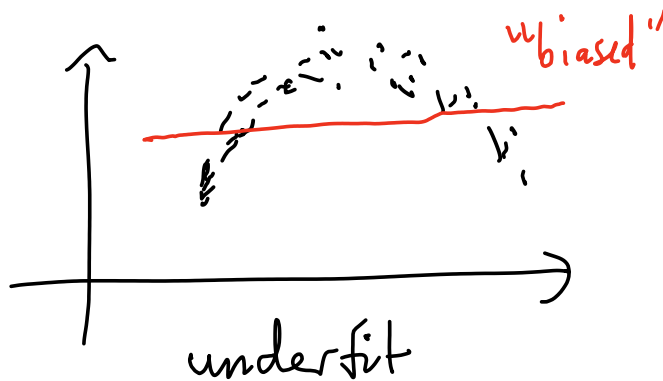
Over fitting



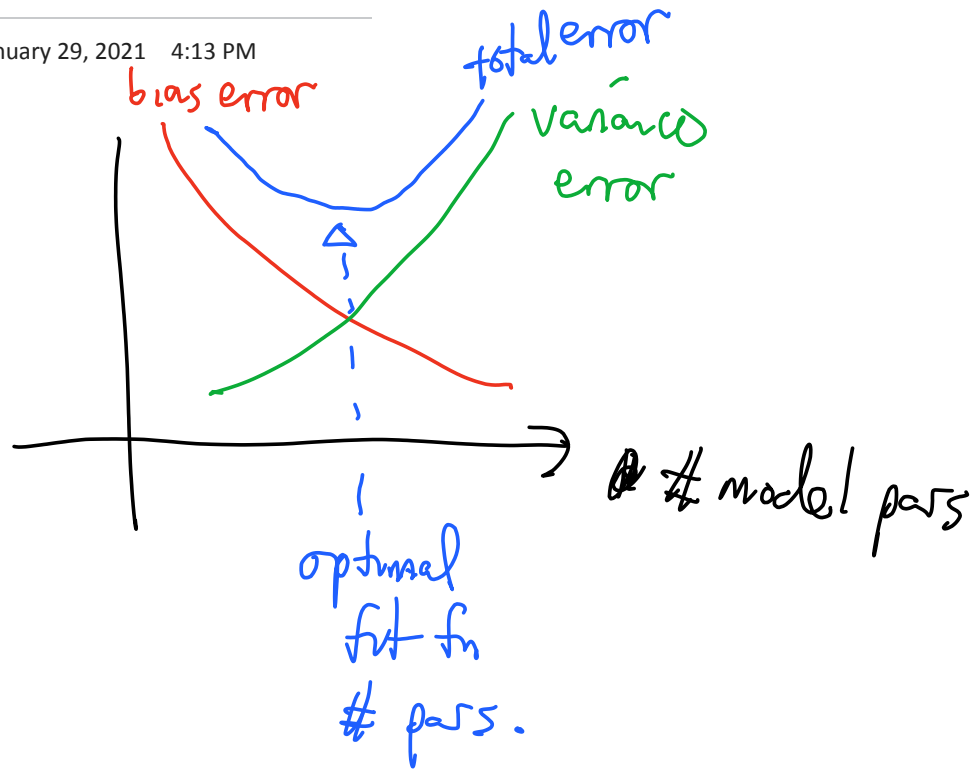
if fit fn has too many parameters ("overparameterized")

→ fit fn can "memorize" training data
→ overfitting

"Bias - variance tradeoff"



Friday, January 29, 2021 4:13 PM

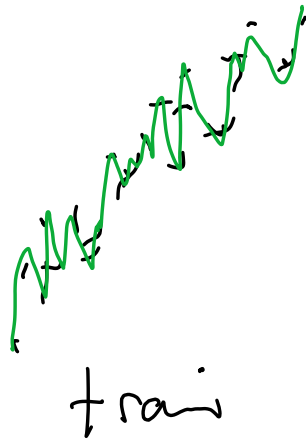


How to diagnose & quantify & prevent overfitting?

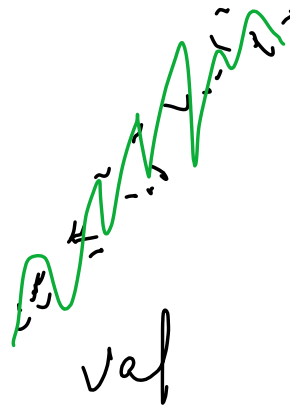
→ validation dataset, (and test dataset)

Friday, January 29, 2021 4:17 PM

val. set: statistically indep. dataset that is not used during training (except to diagnose overfitting)



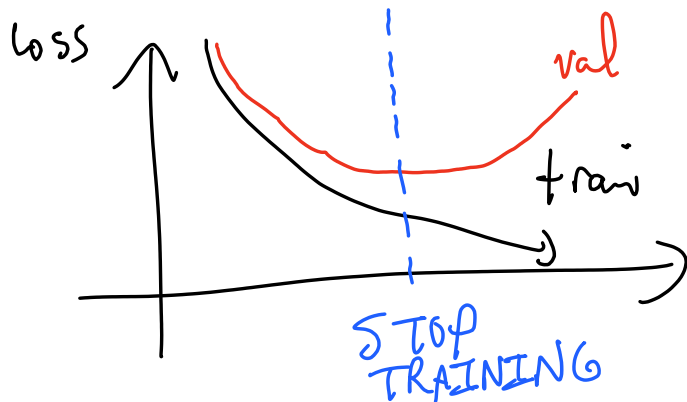
train loss $\rightarrow 0$



val loss $\nrightarrow 0$

may increase w/ time bad

"poor generalization"



epoch "early stopping"

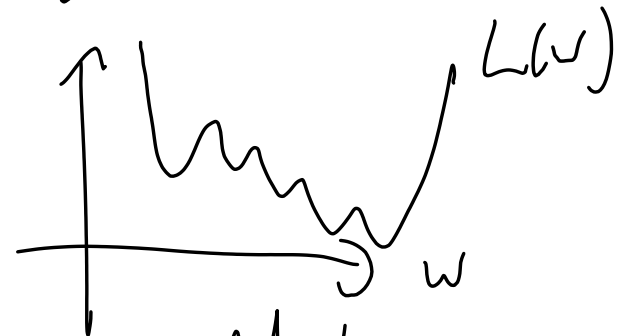
Friday, January 29, 2021 4:24 PM

NN training → (stochastic) gradient descent
 ↘ backpropagation

— What is the goal of training?

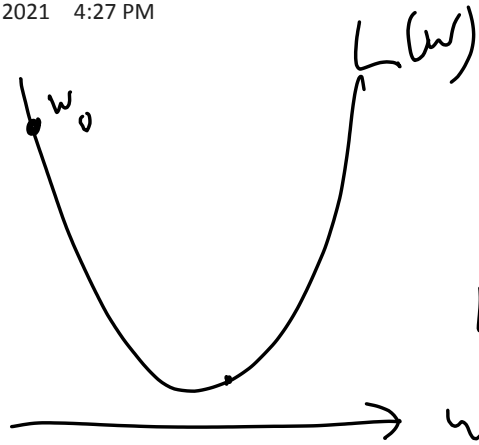
minimize L w.r.t. weights w

wait: $\frac{\partial L}{\partial w} \approx 0$.
 \implies



might be many
 local minima!
How not to get stuck?

Friday, January 29, 2021 4:27 PM



warm up: how to find min. of this fn numerically?

↳ Initial guess

$$2. L(w) \approx L(w_0) + \delta w \left. \frac{\partial L}{\partial w} \right|_0 + \frac{1}{2} \delta w^2 \left. \frac{\partial^2 L}{\partial w^2} \right|_0 + \dots$$

Newton method: 2nd order

$$\delta w \approx - \frac{\left. \frac{\partial L}{\partial w} \right|_0}{\left. \frac{\partial^2 L}{\partial w^2} \right|_0}$$

3. iterate

Problem: doesn't scale well w/ $w \rightarrow \vec{w}$ many parameters
 Because $\frac{\partial^2 L}{\partial w^2}$ becomes matrix that we have to invert.

Friday, January 29, 2021 4:33 PM

So for NNs we make do w/ 1st order method

$$\text{Gradient descent: } \delta w \approx -\alpha \left. \frac{\partial L}{\partial w} \right|_0$$

$$\delta \vec{w} \approx -\alpha \vec{\nabla} L|_0$$

fixed constant parameter
"learning rate"

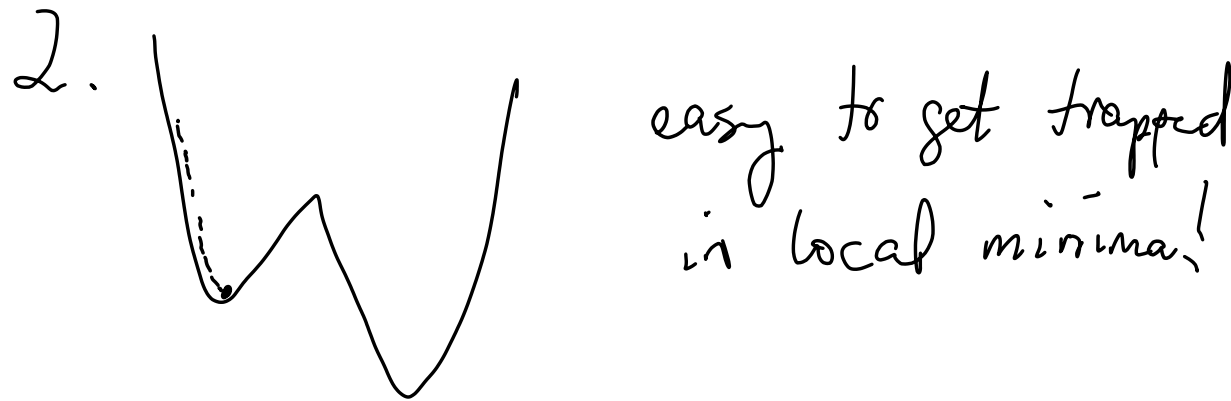


Friday, January 29, 2021 4:38 PM

Problems w/ Vanilla GA:

1. $L = \sum_{x \in \text{Data}} L(x)$

Need to evaluate L on entire data
for a single weight update. Expensive & slow!



easy to get trapped
in local minima!

→ "stochastic
gradient descent"