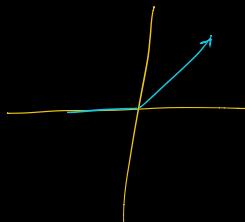


for every weight
we multiply it with
"p" during our time.

↳ dropout is like a layer
in-between. That's how
we put it in keras or TF.

$\Rightarrow \text{ReLU} \rightarrow \text{Rectified Linear Unit.}$

$$(\rho(z) = \max(0, z))$$



$$\frac{\partial \rho(z)}{\partial z} = \begin{cases} 0 & \text{if } z \leq 0 \\ 1 & \text{otherwise.} \end{cases}$$

Even though this is not
differentiable, we don't give a
fuck, we just use a hack
to make it work.

→ ReLU is too good
compared to tanh or σ

↓
It converges a lot
faster than them. *no vanishing
gradient.

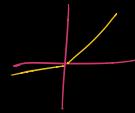


→ for ReLU it's basically
no calculation for getting the f.
It just makes backprop so

that's much faster

$$\because \frac{df}{dw} = 0$$

this doesn't allow gradients to flow through. So everything just piled up.



\Rightarrow That's why Leaky ReLU: $w_i^k = \begin{cases} 0.001z & \text{if } z < 0 \\ z & \text{otherwise} \end{cases}$
(\hookrightarrow doesn't become zero.)

Initialize the weights \rightarrow random distribution.

$w_{ij}^k \neq 0 \rightarrow$ this creates stuck up

① \rightarrow All other activations in further layers will be zero.

② Symmetric. So each weight will have same gradient.

③ init should not have too large or too small. \Rightarrow bad gradients

\rightarrow for ensembles, the more different the base models are, the better the outcome. which is why asymmetry is important with init.

It's also important to normalize the weights itself to prevent exploding / vanishing gradients



works well for σ $\leftarrow w_{ij}^k \sim \text{Uniform}_{\text{Random}} \left[\frac{-1}{\sqrt{f_m}}, \frac{+1}{\sqrt{f_m}} \right]$

this is all empirical, with little theoretical support behind it.

Xavier/Glorot init.

a) Normal $w_{ij}^k \sim N(0, \sigma)$

$$\sigma = \sqrt{\frac{2}{f_m + f_{out}}}$$

b) Uniform

$$w_{ij}^k \sim \text{Uniform} \left[\frac{-\sqrt{6}}{\sqrt{f_m + f_{out}}}, \frac{+\sqrt{6}}{\sqrt{f_m + f_{out}}} \right]$$

Ways You Train

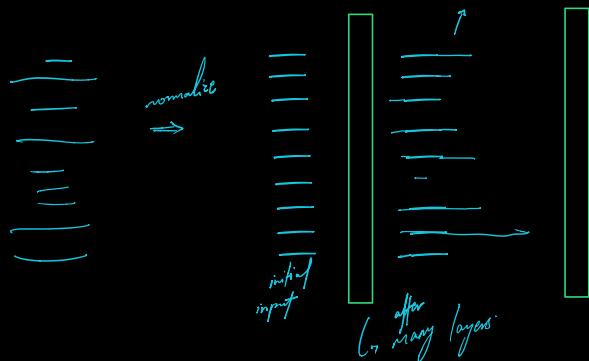
He int \rightarrow works well
with ReLU

(a) Normal $w_j^k \sim N(0, \sigma)$

$$\sigma = \sqrt{\frac{2}{f_m}}$$

(b) Uniform $w_j^k \sim U\left[-\sqrt{\frac{\epsilon}{f_m}}, +\sqrt{\frac{\epsilon}{f_m}}\right]$

Batch Normalization.



\rightarrow Batch normalization means

norm not only with input but

also after every layer.

$$\text{std.dev. } \hat{x}_i = \frac{x_i - \mu \xrightarrow{\text{mean}}}{\sqrt{\sigma^2 + \epsilon}} \xrightarrow{\text{to prevent divide by zero}}$$

$$x_i^k = \gamma \hat{x}_i + \beta$$

\hookrightarrow these two
are learned.

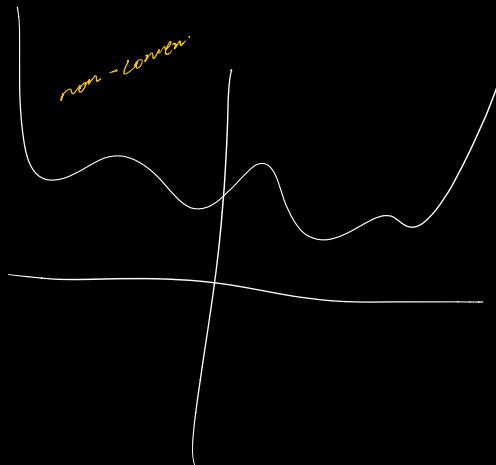
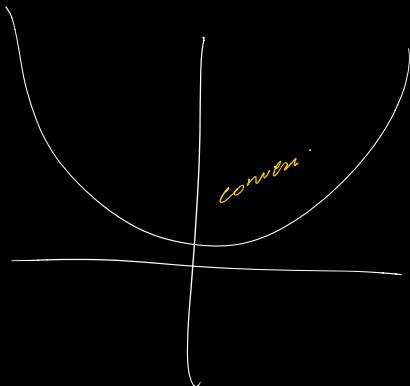
This thing happens for each feature during a batch.

BN also works as a weight regularization.

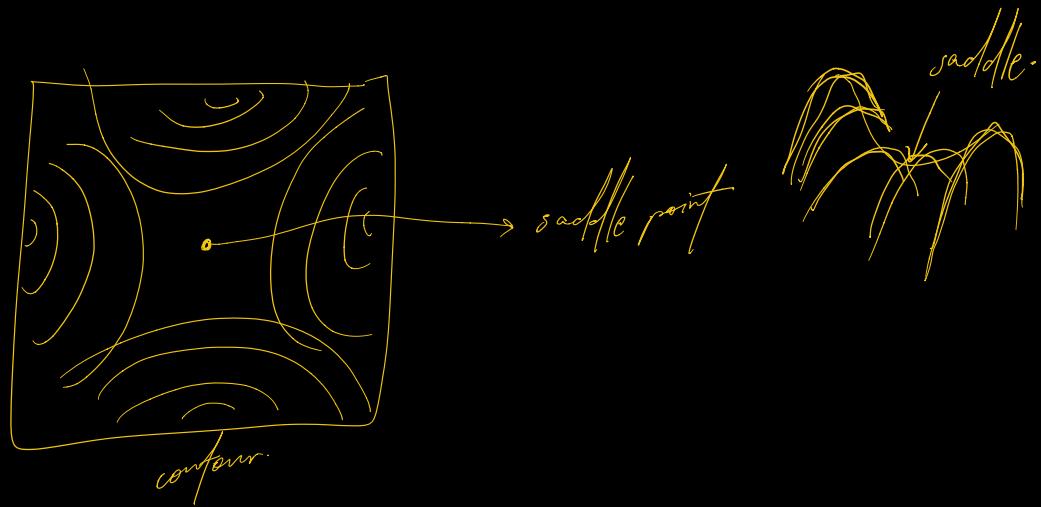
It prevents internal covariate shift and thus allows training deeper neural networks.

✓ # Convex $f^n \rightarrow$ has only one minima.
LR, SVM, etc. local minima = global minima

Non convex $f^n \rightarrow$ multiple local minima.
Deep learning



have to be careful not to settle on a local minima.



Momentum with SGD.

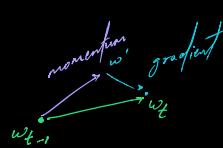
$$\begin{aligned} v_t &= \gamma \hat{v}_t \\ \hat{v}_t &= \gamma \hat{v}_{t-1} + \alpha t \\ &\quad \curvearrowleft 0 < \gamma \leq 1 \\ &\quad \curvearrowleft \text{typically } 0.9 \end{aligned}$$

$$w_t = w_{t-1} - [\gamma v_{t-1} + \eta \hat{g}_t]$$



here we take the addition
of momentum & grad.
both are applied together.

Nesterov Accelerated Gradient (NAG)



Here we apply momentum first,
then compute grad at that point(w')
then apply it.

$$w_t = w_{t-1} - (\gamma v_{t-1} + \eta g')$$

$$g' = \left(\frac{\partial L}{\partial w} \right)_{w'}$$

$$w' = w_t - \gamma v_{t-1}$$

momentum applied first
& then grad

* NAG has kind the same performance as SGD + momentum. Nothing fancy.

Ada. Grad.

→ Adaptive learning rate

$$w_t = w_{t-1} - \eta'_t \hat{g}_t$$

$\eta'_t = \frac{\eta}{\sqrt{d_{t-1}} + \epsilon}$ even η is time dependent now.

just to avoid divide by zero

$$d_{t+1} = \sum_{i=1}^t \hat{g}_i^2 \quad \left. \begin{array}{l} \text{all grads up to} \\ \text{this point.} \end{array} \right\} \text{this keeps on } \uparrow \quad \text{and thus } \eta' \text{ reduces with time.}$$

↓
always positive

⊕ → No need to fiddle with η

⊕ → sparse features benefit from it

⊖ → η does ↓ but then later on it gets super low as $t \uparrow$

↓ slow convergence

AdaDelta & RMSProp

$$w_t = w_{t-1} - \eta'_t \hat{g}_t$$

$$\eta'_t = \frac{\eta}{\sqrt{d_{t-1}} + \epsilon}$$

$$d_{t-1} = \left[\hat{g}_{t-1} \hat{g}_{t-1}^T + (1-\gamma) d_{t-2} \right]$$

normally 0.9

$$d_{t+1} = 0.1 \hat{g}_{t-1}^2 + 0.9 \left[0.1 \hat{g}_{t-2}^2 + 0.9 [d_{t-3}] \right] \quad \text{so on...}$$

→ was controlling the growth of denominator

→ avg. weighted avg. of \hat{g}_i^2
→ instead of own \hat{g}_i^2

→ The updates still remain from AdaGrad

Adam
[Adaptive Moment Estimation]

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \quad 0 \leq \beta_1, \beta_2 \leq 1$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\hat{m}_t = \frac{m_t}{1 - (\beta_1)^t} \quad \hat{v}_t = \frac{v_t}{1 - (\beta_2)^t}$$

$$w_t = w_{t-1} - \alpha \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}$$

Monitoring gradients. \rightarrow should do it for each layer, epoch \rightarrow good practice

Vanishing/Exploding gradients
can be affected.

\hookrightarrow can be avoided using grad clipping. $\rightarrow \alpha_{new} = \left(\frac{\alpha}{\|g\|_2} \right) * \gamma$ threshold
 $\hookrightarrow \gamma < 1 * \gamma$

All libraries have
this L2 norm clipping
feature.

Softmax

o only works for binary classification.

we want the prob to add up to 1
we want multiplex

This is just an extension of LR
↳ generalization

$$\sigma(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{e^z + 1}$$

$$= \underbrace{e^z}_{=1} \quad \left| \begin{array}{l} \\ \end{array} \right. \quad \sum_j \epsilon_j$$

$$-\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^k y_{ij} \log p_{ij}$$

$$\sigma_1(z_1) = \frac{e^{z_1}}{\sum_{i=1}^k e^{z_i}}$$

$$\sigma_2(z_2) = \frac{e^{z_2}}{\sum_{i=1}^k e^{z_i}}$$

$$\sigma_3(z_3) = \frac{e^{z_3}}{\sum_{i=1}^k e^{z_i}}$$

For every point
& for every class we get loss

→ This is an extension of
↳ softmax

This is for classification,

for regression we make it
a linear wif

$$so f(x) = z = w^T x$$

And loss is squared loss

#How to Train an MLP (summary)

- ① Preprocess \rightarrow Normalise
- ② weight init \rightarrow Xavier/Glorot \rightarrow Sigmoid/Tanh
 \rightarrow He \rightarrow ReLU
 \rightarrow Gaussian with small σ
- ③ Activation \rightarrow ReLU
- ④ BN and Dropout
- ⑤ Optimizer \rightarrow Adam $\#$ best
- ⑥ Hyperparameters \rightarrow Hayes
 $\#$ neurons
 $\#$ dropout rate
- ⑦ Loss \rightarrow squared loss \rightarrow regression.
 \rightarrow multi-class log loss \rightarrow classification
- ⑧ Monitor gradients \rightarrow Clipping
- ⑨ Plot using Tensorboard \rightarrow 
- ⑩ Avoid overfitting.

Even among one of these points is enough to throw us off

#Denoising Autoencoder

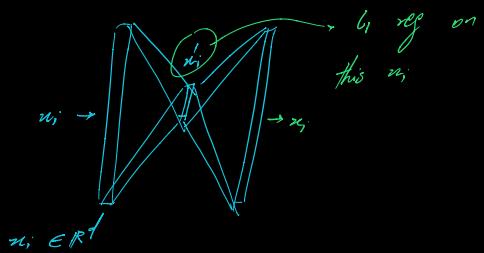
$$\tilde{x}_i = x_i + N(0, \sigma)$$

$$x_i^{\text{corrupted}} \Rightarrow \boxed{\text{[Input Image]}} \Rightarrow x_i$$

Data is intentionally corrupted
 and AE is trained to get pure data
 from slightly corrupted data.

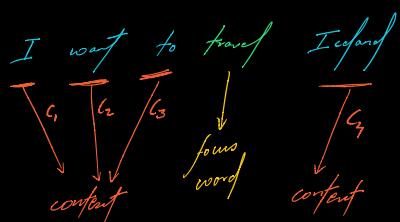
Sparse AF

If feature vector is too sparse



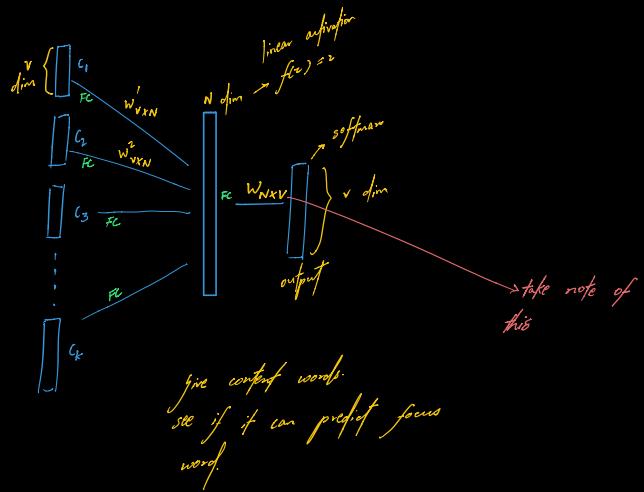
If linear activations are used and
only one sigmoid hidden layer is
taken, the result is kinda similar
to what PCA would give us.

WordVec \rightarrow not a deep learning
algo



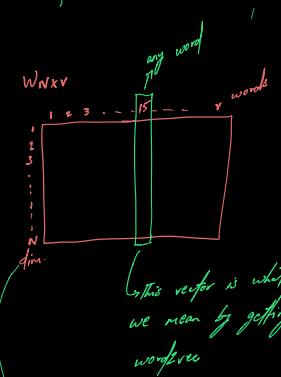
(CBOW)
Collective Bag of Words.

- ① $v = \text{size of vocab}$
- ② one hot encode each word
 $w_i \in \mathbb{R}^v$

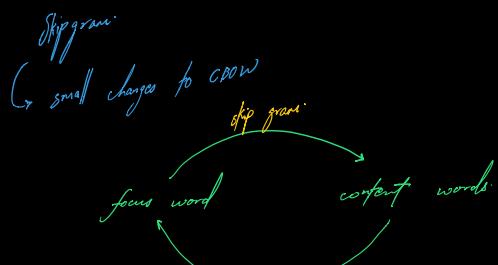


Training data:

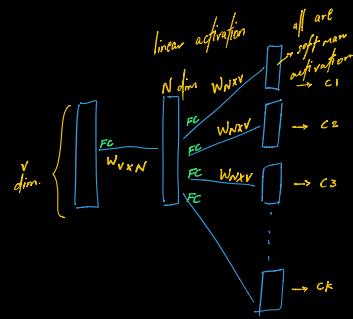
focus	content
framed	{ I, want, to, Isolate }



Dimension of hidden layer. Can be changed depending on need.
Also the dimension of wordvector.



CBOW



The no. of weights are same
with CBOW & skipgram. $(k+1)(N \times V)$

kinda easy:
content given \rightarrow predict one word

hard one word given \rightarrow predict content

\hookrightarrow more powerful

CBOW \rightarrow 1 softmax \rightarrow fast \rightarrow good for freq.
words. cor well
have more context words

Skipgram \rightarrow k softmax \rightarrow slow
 \downarrow
words with less data and

$k \rightarrow$ # context words.

$k \uparrow \rightarrow$ more context \rightarrow better result.

$N \rightarrow$ # dim.

$N \uparrow \rightarrow$ more information.

$$\# \text{ weights} = (k+1)(N \times V)$$

$$N = 200 \text{ dim}$$

$$k = 5$$

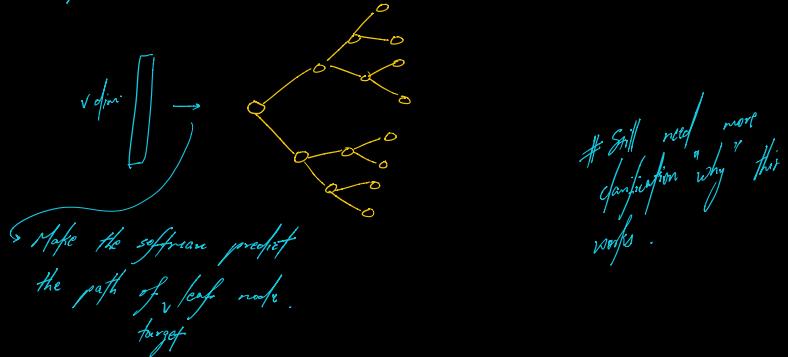
$$V = 10k \text{ words}$$

$$\text{Still } \# \text{ weights} = 12 \text{ million.}$$

need to optimize!!

Hierarchical
softmax

optimize V-softmax somehow.



If still not clear why this works.

So we converted $O(v)$ to $O(\log v)$

Optimization 2

Negative Sampling



→ always update this one.

for the others we have
a probability that they will be
updated.

$$P(w_i) = 1 - \frac{c}{\sqrt{\text{freq.}(w_i)}} \quad \text{generally } 10^{-5}$$

This is only intuition for hierarchical softmax and negative sampling.

/ / / / /