



UAI
UNIVERSIDAD ADOLFO IBÁÑEZ

Extras:

Bash Normalization, Dropout y Cross Entropy

Aprendizaje profundo

Dr. Juan Bekios Calfa

Magister en *Data Science*
Facultad de Ingeniería y Ciencias

Información de Contacto

- Juan Bekios Calfa
 - email: juan.bekios@edu.uai.cl
 - Web page: <http://jbekios.ucn.cl>
 - Teléfono: 235(5162) - 235(5125)



Contenidos

Introducción

Batch normalization

Definición

Pytorch

Regularización

Definición

Dropout

Pytorch

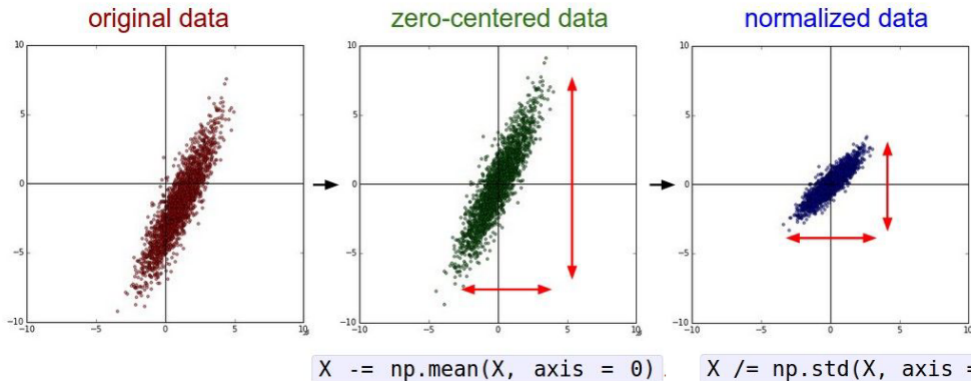
Clasificación multiclase

Tipos de clasificación

Funciones de activación: Softmax

Funciones de pérdida

Preprocesamiento de datos



(Asuma que $X [N \times D]$ es una matriz de datos, cada fila es un ejemplo)

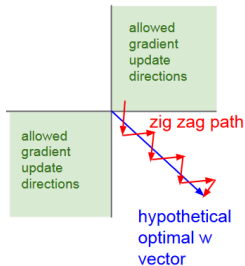
Preprocesamiento de datos

Consideremos lo que sucede cuando la entrada a una neurona es siempre positiva

$$f\left(\sum_i w_i x_i + b\right)$$

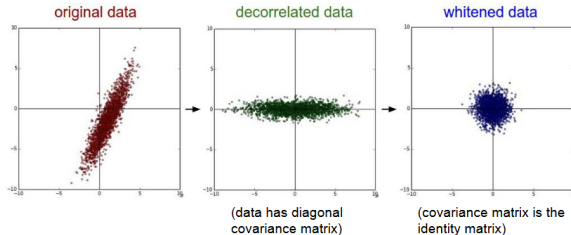
¿Que podemos decir de los gradientes en \mathbf{w} ?

¿Siempre positivos o negativos?



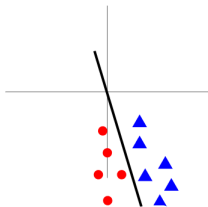
Preprocesamiento de datos

En la práctica, es posible ver el uso de ACP (**PCA**) y blanqueamiento (**Whitening**) de los datos.

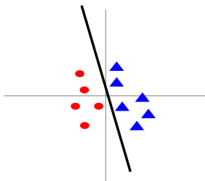


Preprocesamiento de datos

Antes de la normalización: pérdida de clasificación muy sensible a cambios en la matriz de peso, difícil de optimizar



Después de la normalización: menos sensible a pequeños cambios de peso, más fácil de optimizar



Resumen: En la práctica para imágenes

Por ejemplo: Consideremos CIFAR-10 con imágenes de $[32,32,3]$

- Restar la imagen promedio (arreglo de $[32,32,3]$)
- Restar el promedio por canal (3 números)
- Restar el promedio por canal y dividirlo por la desviación estándar por canal (3 números)

Batch normalization

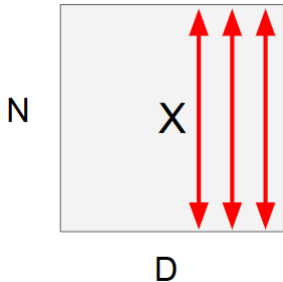
¿Quieres activaciones de varianza unitaria de media cero?

Consideremos un *batch* de activaciones en alguna capa. Para hacer que cada dimensión varíe la unidad de media cero, aplique:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Batch normalization

Input: $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} \text{ Promedio por canal, } [D]$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 \text{ Varianza por canal } [D]$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}} \text{ x normalizado, } [N,D]$$

¿Qué pasaría si el uso de media cero y variancia unitaria es una restricción demasiada fuerte?

Batch normalization

Input: $x : N \times D$ Parámetros de cambio y escala que se pueden aprender: $\gamma, \beta : D$
Aprendiendo los parámetros
 $\gamma = \sigma, \beta = \mu$ permitirá recuperar la

función identidad

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} \text{ Promedio por canal, [D]}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 \text{ Var por canal [D]}$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}} \text{ x normalizado, [N,D]}$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j \text{ Salida, [N,D]}$$

Batch normalization - durante evaluación

Input: $x : N \times D$ Parámetros de cambio y escala que se pueden aprender: $\gamma, \beta : D$
Durante la evaluación de la red *batchnorm* se transforma en un operador lineal $\mu_j =$ Promedio de los

valores vistos durante el entrenamiento Promedio por canal, $[D]$
 $\sigma_j^2 =$ Promedio de los valores vistos durante el entrenamiento Var por canal $[D]$
 $\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$ x normalizado, $[N, D]$
 $y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$ Salida, $[N, D]$



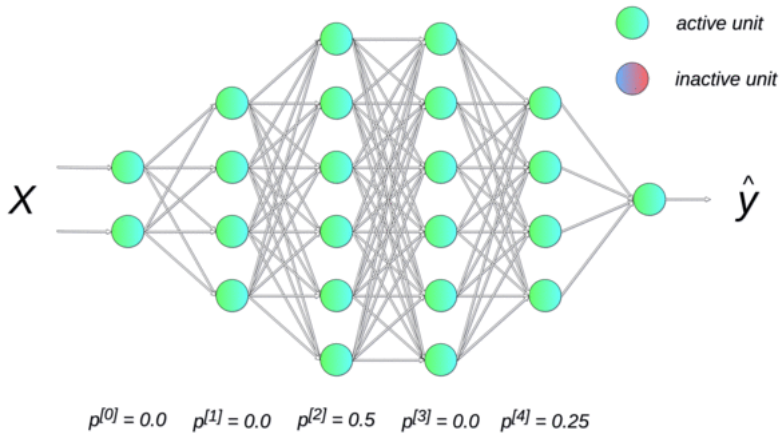
Batch Normalization: Pytorch

```
class Multilayerperceptron(nn.Module):  
  
    def __init__(self):  
        super().__init__()  
        self.layers = nn.Sequential(  
            nn.Flatten(),  
            nn.Linear(32 * 32 * 3, 64),  
            nn.BatchNorm1d(64),  
            nn.ReLU(),  
            nn.Linear(64, 32),  
            nn.BatchNorm1d(32),  
            nn.ReLU(),  
            nn.Linear(32, 10)  
        )
```

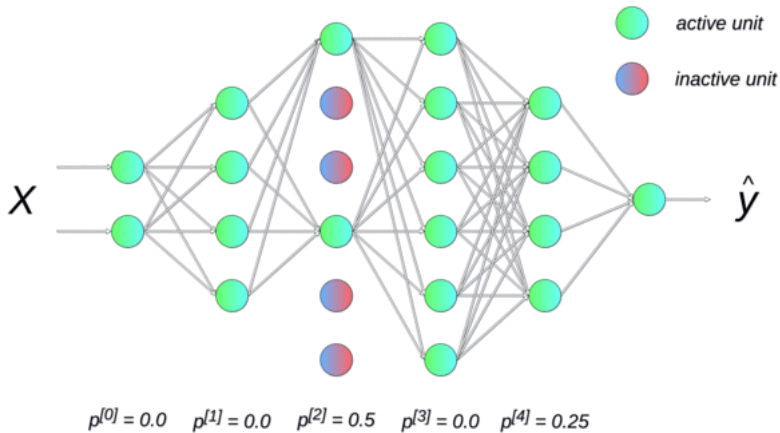
Dropout

Dropout es un algoritmo para entrenar redes neuronales mediante la eliminación aleatoria de unidades durante el entrenamiento para evitar su coadaptación.

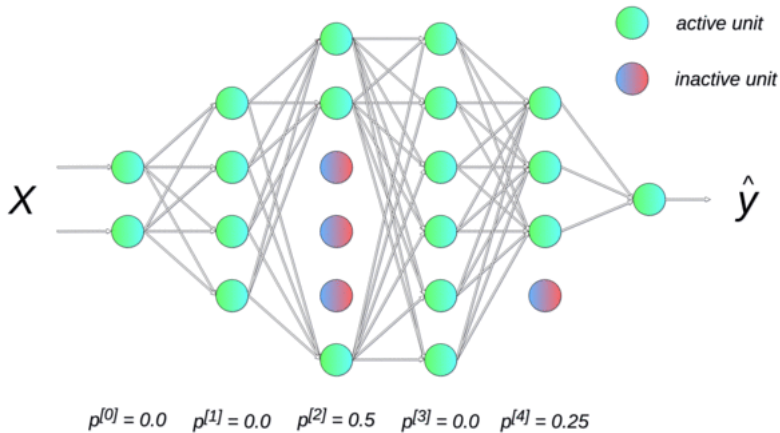
Dropout (I)



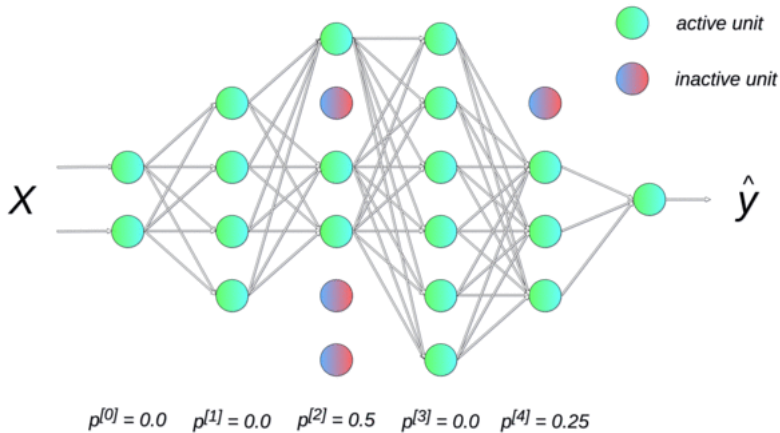
Dropout (II)



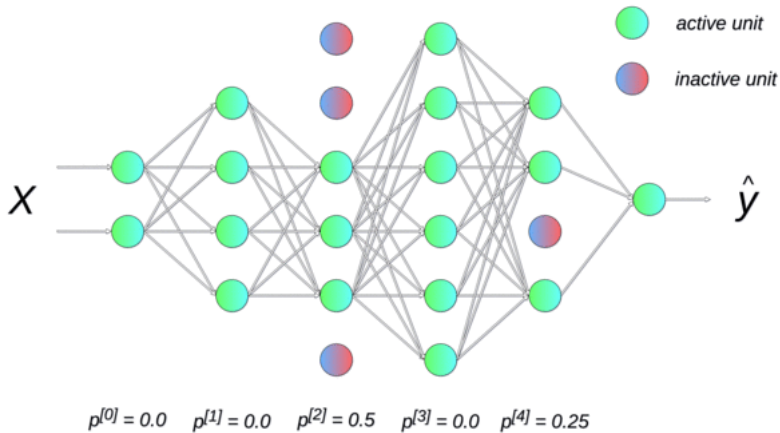
Dropout (III)



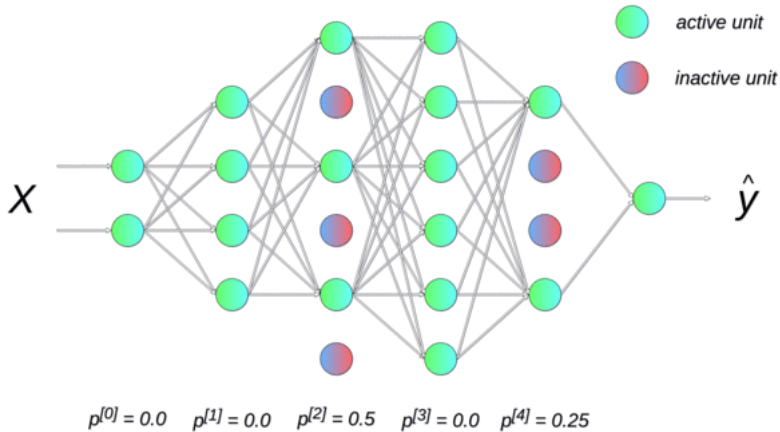
Dropout (IV)



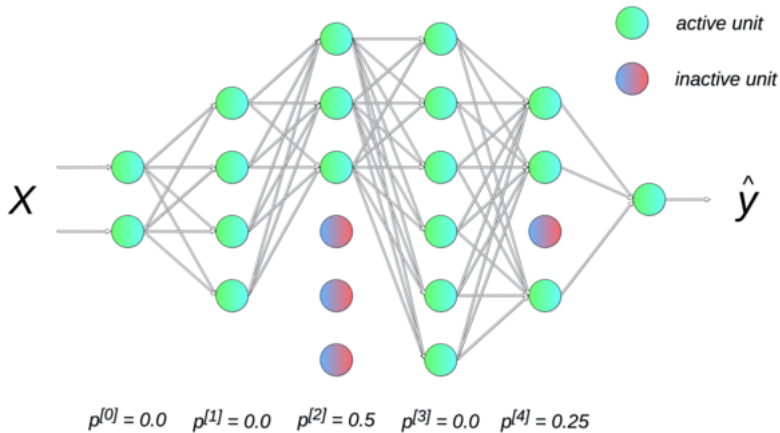
Dropout (V)



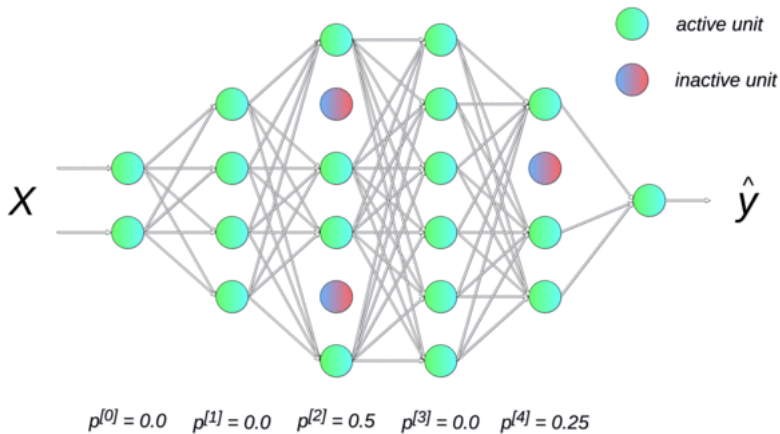
Dropout (VI)



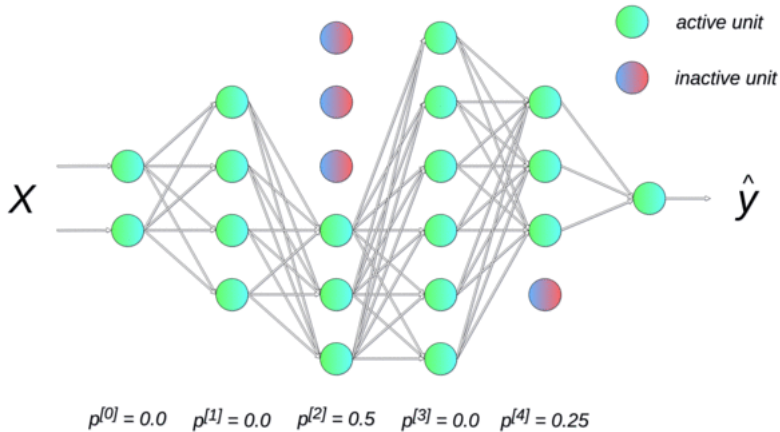
Dropout (VII)



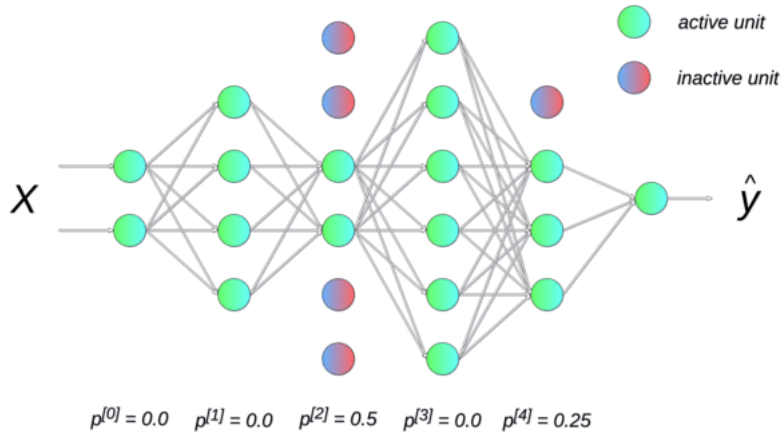
Dropout (VIII)



Dropout (IX)



Dropout (X)



Dropout: Pytorch

Regularization 1: Dropout

Dropout: in pytorch is implemented as `torch.nn.Dropout`

If we have a network:

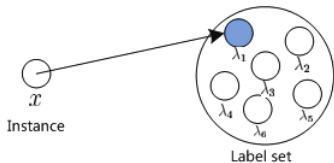
```
model = torch.nn.Sequential(  
    torch.nn.Linear(1,100), torch.nn.ReLU(),  
    torch.nn.Linear(100,50), torch.nn.ReLU(),  
    torch.nn.Linear(50,2))
```

We can simply add dropout layers:

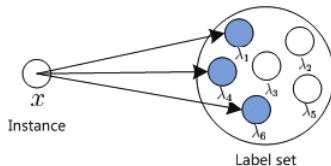
```
model = torch.nn.Sequential(  
    torch.nn.Linear(1,100), torch.nn.ReLU(),  
    torch.nn.Dropout(),  
    torch.nn.Linear(100,50), torch.nn.ReLU(),  
    torch.nn.Dropout(),  
    torch.nn.Linear(50,2))
```

Note: A model using dropout has to be set in **train** or **eval** model.

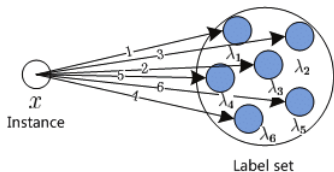
Clasificación multiclase



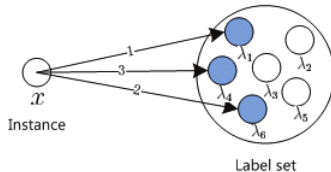
(a) multiclass classification



(b) multilabel classification



(c) label ranking



(d) multilabel ranking

Función de activación: *Softmax*

Input pixels, x



Shape: (3, 32, 32)

Feedforward output, y_i

	cat	dog	horse
Forward propagation	5	4	2
	4	2	8
	4	4	1

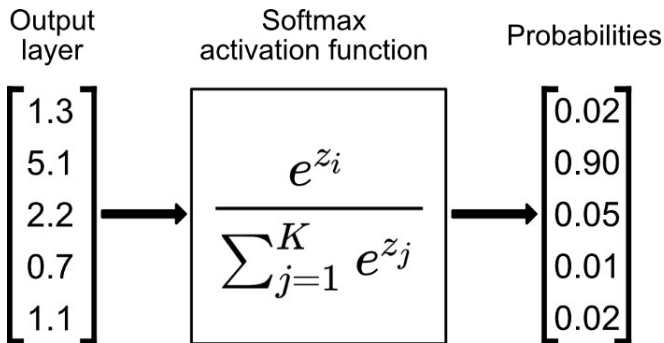
Shape: (3,)

Softmax output, $S(y_i)$

	cat	dog	horse
Softmax function	0.71	0.26	0.04
	0.02	0.00	0.98
	0.49	0.49	0.02

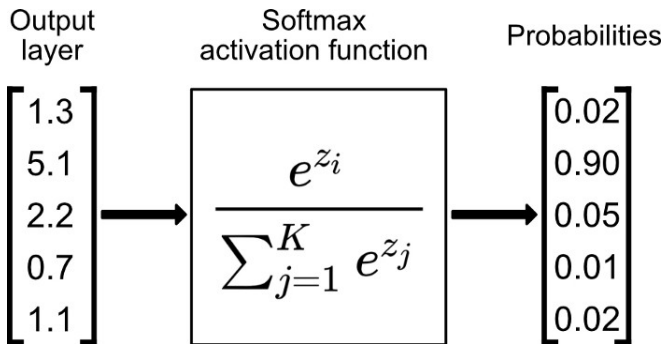
Shape: (3,)

Función de activación: *Softmax*



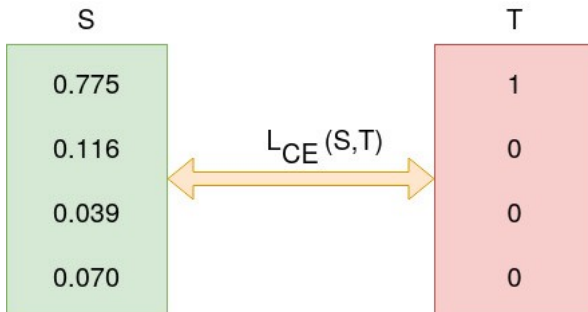
<https://towardsdatascience.com/softmax-activation-function-explained-a7e1bc3ad60>

Función de pérdida: *Cross Entropy*



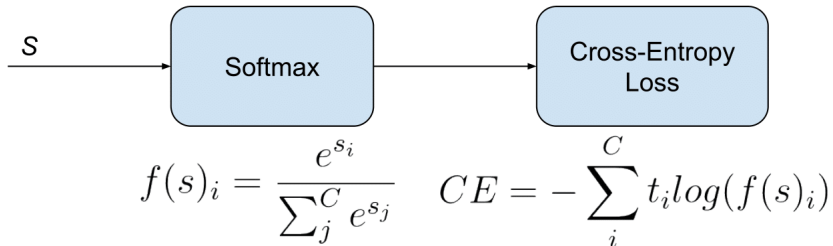
<https://towardsdatascience.com/softmax-activation-function-explained-a7e1bc3ad60>

Función de pérdida: *Cross Entropy*



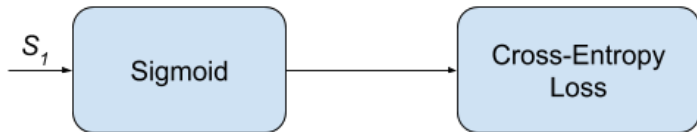
<https://towardsdatascience.com/cross-entropy-loss-function-f38c4ec8643e>

Función de pérdida: *Cross Entropy*



https://gombru.github.io/2018/05/23/cross_entropy_loss/

Función de pérdida: *Binary Cross-Entropy*



$$f(s_i) = \frac{1}{1 + e^{-s_i}}$$
$$CE = -t_1 \log(f(s_1)) - (1 - t_1) \log(1 - f(s_1))$$

https://gombru.github.io/2018/05/23/cross_entropy_loss/

¿Preguntas?