

Capítulo 4 : Neurônios

Em março de 2016, o software AlphaGo venceu um mestre de Go. O feito é impressionante por se tratar de um jogo difícil de computar.

Inventado há mais de 2,500 anos, motivou avanços em matemática. Existem $2,08 * 10^{170}$ maneiras válidas de dispor as peças no tabuleiro. O polímata chinês Shen Kuo (1031–1095) chegou a um resultado próximo 10^{172} séculos atrás. Vale lembrar que o número de átomos no universo observável é de módicos 10^{80} .

No capítulo anterior, aprendemos uma formulação básica de modelo preditivo, a regressão linear simples. A seguir, estenderemos nosso leque de ferramentas para novas classes de relações, também incluindo mais informações na entrada de nossos modelos.

Mais do que isso, conheceremos a primeira máquina inteligente da história.

O perceptron de Rosenblatt

Frank Rosenblatt (1928 - 1971) nasceu e morreu em 11 de julho, mas esse não é o fato mais curioso da biografia deste psicólogo. Foi o responsável pelo desenvolvimento do primeiro neurônio artificial. Em suas palavras, o primeiro objeto não biológico a recriar uma organização do ambiente externo com significado.

It can tell the difference between a cat and a dog, although it wouldn't be able to tell whether the dog was to the left or right of the cat. Right now it is of no practical use, Dr. Rosenblatt conceded, but he said that one day it might be useful to send one into outer space to take in impressions for us. - New Yorker, December, 1958¹

O aparato reproduzia o entendimento da época sobre o funcionamento de um neurônio. O corpo recebe sinais de dendritos e, após processamentos ocultos, produz um output na forma de sinal elétrico pelo axônio. A primeira matematização viria do modelo de McCulloch & Pitts (“A Logical Calculus of the Ideas Immanent in Nervous Activity”, 1943).

Em 1949, Donald Hebb descreveu em seu clássico *The Organization of Behavior* um mecanismo plausível para a aprendizagem. Comumente expressa na máxima “Cells that fire together wire together” (células que disparam juntas, conectam-se entre si).

Com o objetivo de criar uma máquina que pudesse processar inputs diretamente do ambiente físico (e.g. luz e som), Rosenblatt concebeu uma extensão elegante do modelo em 1957 (“The Perceptron[do latim, *percipio, compreender*] – a perceiving and recognizing automaton. Report 85-460-1, Cornell Aeronautical Laboratory”). Composto de três partes: o sistema S (sensório); o sistema A (associação) e o sistema R (resposta). O neurônio “lógico” criado por McCulloch & Pitts foi modificado de maneira a processar inputs através de pesos antes da saída. A aprendizagem se dá pela modificação desses pesos.

Inicialmente, o perceptron foi simulado em um IBM 704 (também berço das linguagens FORTRAN e LISP). Em seguida, implementado como um dispositivo físico, batizado de Mark I Perceptron.² Um estudo mais profundo foi publicado por ele em 1962 (*Principles of neurodynamics*).

¹Ele consegue diferenciar um gato de um cachorro, ainda que não seja capaz de dizer se o cachorro estava à esquerda ou à direita do gato. No momento, não tem uso prático, Dr. Rosenblatt admitiu, porém disse que um dia pode ser útil para enviar um [aparato] ao espaço para capturar impressões para nós.

²Mark I é um título comumente utilizado para a primeira versão de uma máquina.

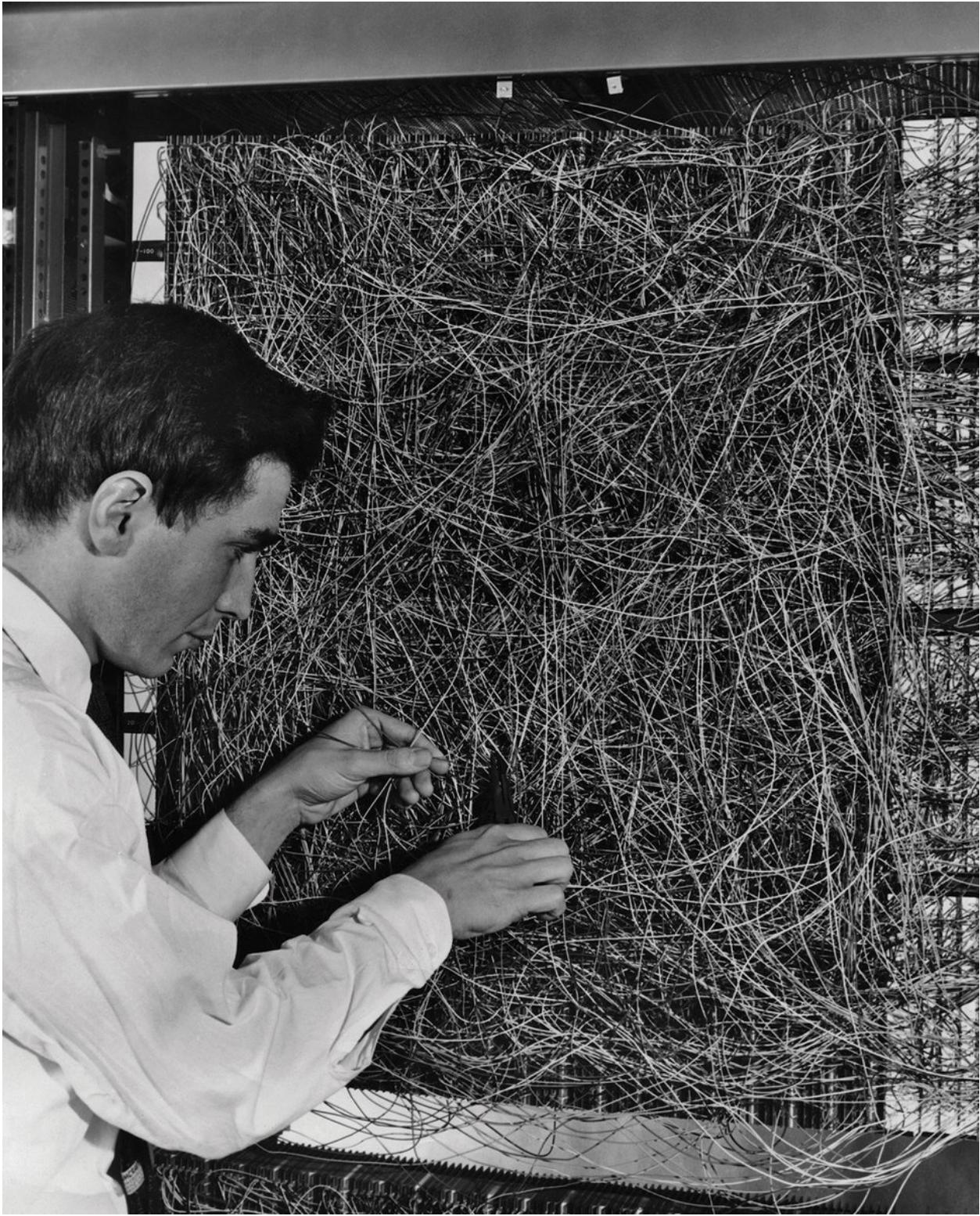


Figure 1: Frank Rosenblatt e Mark I.

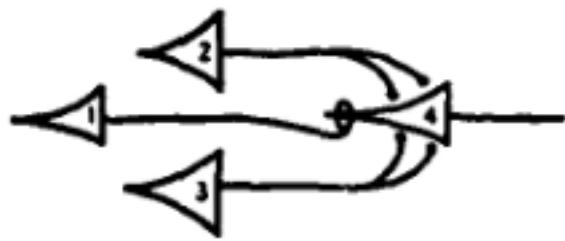


Figure 2: Diagrama de células lógicas em McCulloch & Pitts

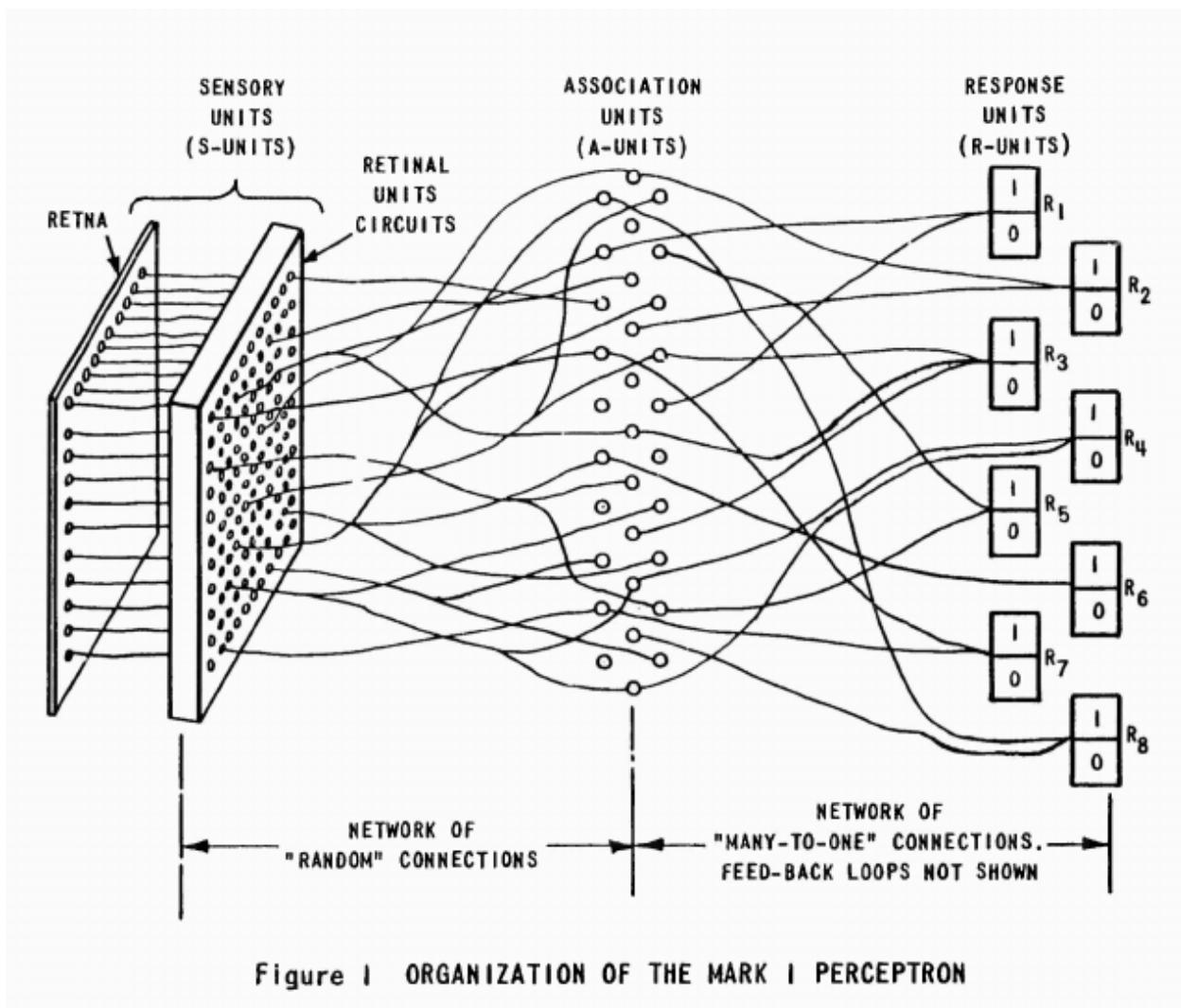


Figure 3: Organização do Mark I, retirado de seu manual de uso original

Rosenblatt protagonizava calorosos debates sobre inteligência artificial na comunidade científica junto a Marvin Minsky, um amigo da adolescência. Em 1969, Minsky e um matemático (Seymour Papert) publicaram um livro centrado no Perceptron (Perceptrons: An Introduction to Computational Geometry). Nele, provaram que o neurônio artificial era incapaz de resolver problemas não-lineares do tipo XOR. Para um problema eXclusive OR (OU eXclusivo) o neurônio deve disparar diante do estímulo A ou do estímulo B, porém não diante de ambos.

O impacto foi devastador sobre o otimismo vigente e se passou um período de 10 anos de baixíssima produção, conhecido como ‘idade das trevas’ do conexionismo. A retomada dos neurônios artificiais aconteceu somente na década de 80. Infelizmente, Rosenblatt morreu prematuramente em 1972 num acidente de barco, não presenciando o renascimento dos perceptrons.

Sabendo das origens do modelo, é curioso que a maioria dos cursos introduzam perceptrons do ponto de vista puramente matemático, apontando a semelhança com neurônios como mera curiosidade. Pelo contrário, a inspiração em neurônios biológicos e posterior sucesso nas tarefas designadas fala em favor de um fantástico caso de sucesso via engenharia reversa.

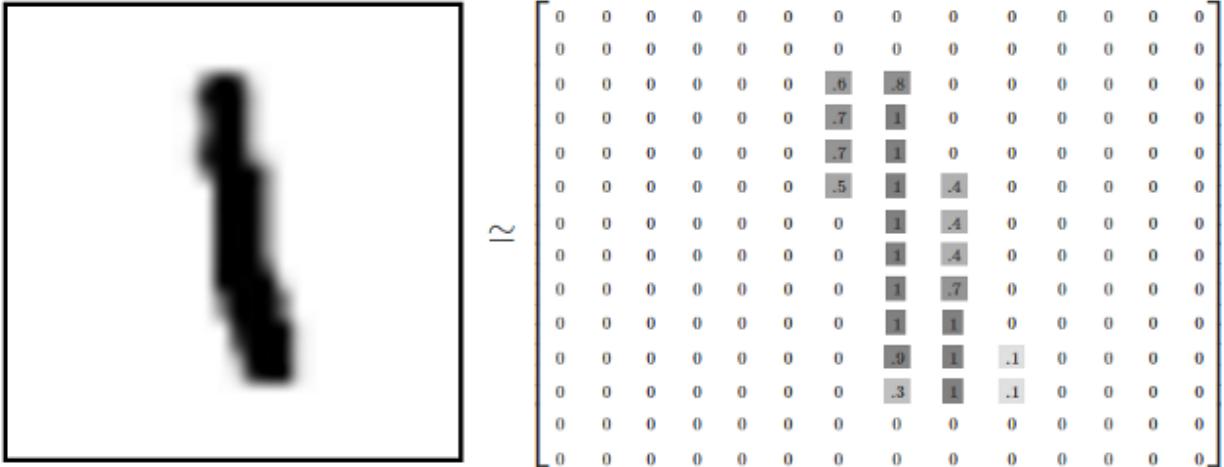
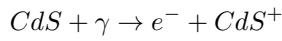


Figure 4: Exemplo de “1” em letra cursiva e sua representação numa matriz 2x2. <http://colah.github.io/posts/2014-10-Visualizing-MNIST/>

Criando neurônios

Mark I foi criado para reconhecimento visual, podendo ser considerado avô da visão computacional. Possuía um campo de entrada fotossensível de 20x20 (400) células de Sulfeto de Cádmio, as unidades S. Ao reagir com a luz, CdS emite um elétron:



Caso a célula seja ativada, envia o sinal eletrônico a uma unidade intermediária A. A unidade intermediária, por sua vez, transmite um sinal eletrônico à saída. **A intensidade do sinal é regulada por sucessos prévios** de maneira a ajustar o sinal para a classificação correta. O aparato físico mimetiza o modelo matemático do **classificador**.

Um sinal luminoso excita cada campo de maneira diferente, ativando células de acordo com a quantidade de luz captada. Matematicamente, representamos cada neurônio sensível à luz como uma célula na matriz de entrada.

O dígito acima ('1') está numa imagem com 14 x 14 pixels (196 valores entre: 1, preto; e 0, branco). Esses pixels podem ser esticados e vistos como uma matriz X de dimensão [196x1] com valores entre 0 e 1 em cada elemento.

Vamos simular uma imagem semelhante:

Um sinal luminoso excita cada campo de maneira diferente, ativando células de acordo com a quantidade de luz captada.

```
>library(magrittr)
>set.seed(2600)
>my.image.data <- c(0,0,0,0,0,0,0,0,0,0,0,0,
+ 0,0,0,0,1,.9,.6,1,0,0,0,0,0,
+ 0,0,0,0,1,0,1,0,0,0,0,0,0,
+ 0,0,0,0,0.9,0,1,1,0,0,0,0,0,
+ 0,0,0,0,0,0,1,1,0,0,0,0,0,
+ 0,0,0,0,0,0,1,1,0,0,0,0,0,
+ 0,0,0,0,0,0,1,1,0,0,0,0,0,
+ 0,0,0,0,0,0,1,1,0,0,0,0,0,
+ 0,0,0,0,0,0,1,1,0,0,0,0,0,
```

```

,0,0,0,0,0,0,1,.9,0,0,0,0,0,0,
,0,0,0,0,0,0,0,0,0,0,0,0,0,0,
,0,0,0,0,0,0,0,0,0,0,0,0,0,0,
,0,0,0,0,0,0,0,0,0,0,0,0,0,0) %>%
matrix(.,14,14,byrow=T)
> image(t(my.image.data[14:1,]), axes = FALSE, col = grey(seq(1, 0, length = 256)))

```



Eis a nossa imagem [14x14]. O computador lê os valores entre 0 (branco) e 1 (branco), dispondo para nós o sinal visual correspondente numa paleta de cores. Aqui usamos 256 tons cinza.

Em regressão linear múltipla, calculamos um peso β para cada variável. O racional aqui é parecido: ponderamos cada pixel por seus respectivos pesos. Em analogia, cada imagem é uma observação de 196 variáveis.

Classificação

Na tarefa de regressão linear, o output deveria ser um número real $Y \beta * X$ com $X, Y \in \mathbb{R}$. Usaremos o perceptron para classificação: as possibilidades de saída são categorias. Isto é, o output é discretizado, geralmente num conjunto binário (e.g. $\{-1, 1\}$ ou $\{0, 1\}$). O neurônio deve disparar (output $y = 1$) caso seja um navio ou permanecer em repouso ($y = -1$) caso não seja.

Matematicamente, é uma multiplicação da matriz de valores da imagem x_j , de dimensão $[100 \times 1]$ por uma matriz $W_{[100 \times 1]}$ que traz os pesos (weights) estimados para cada pixel para cada classe. Então, forçamos o resultado para +1 ou -1 com uma função de ativação (ϕ).

$$y = \phi(W^T X)$$

Usaremos a função *Heaviside step*:

$$\phi(x) = \begin{cases} +1 & se \quad x \geq 0 \\ -1 & se \quad x < 0 \end{cases}$$

Em R:

```

library(magrittr)
# Heaviside
>phi_heavi <- function(x){ifelse(x >=0,1,-1)}
# Iniciando pesos com base em distribuição normal
>my_weights <- rnorm(100)/100

```

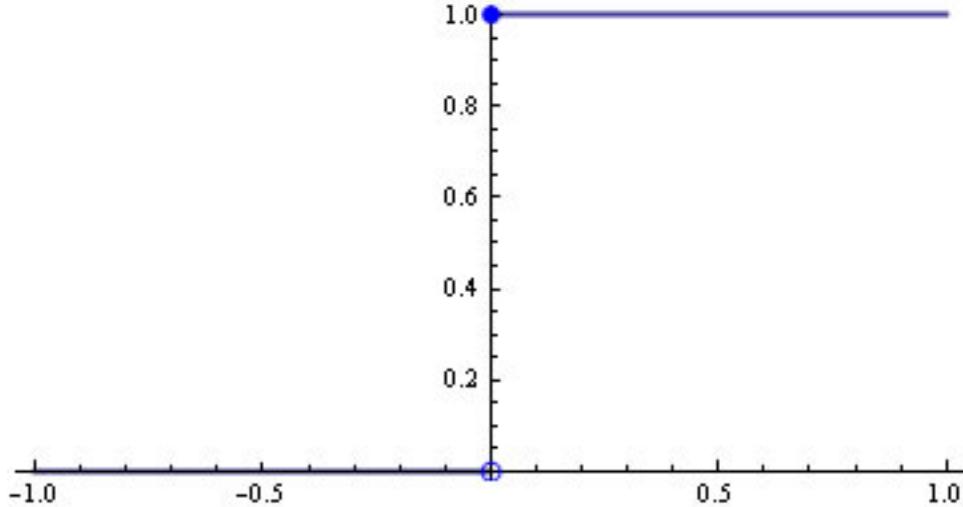


Figure 5: Heaviside step function

```

>w <- matrix(my_weights,100,1)
# Multiplicacao usando o operador %*%
>as.vector(x) %*% w
# Score
[,1]
[1,] 20.19958
# Funcao de ativacao usando %>% para encadeamento
>as.vector(x) %*% w %>% phi_heavi
[,1]
[1,] 1

```

Para o exemplo acima, nosso neurônio com pesos aleatórios foi ativado para o estímulo aleatório x . Inicialmente, estabelecemos pesos aleatórios a partir de uma distribuição normal.

Então, o objetivo é observar as respostas corretas em várias imagens x_i e alterar os valores de W para que os scores maiores sejam os das classes corretas.

O processo de treino é bastante simples:

Seja x_{i_j} o i -ésimo pixel da observação j . E w_0 o peso correspondente inicial, o peso atualizado, w' é:

$$w' = w_0 + \Delta w$$

Em que Δw indica o magnitude e o sentido da modicação no peso.

Aceitemos, por enquanto, a fórmula:

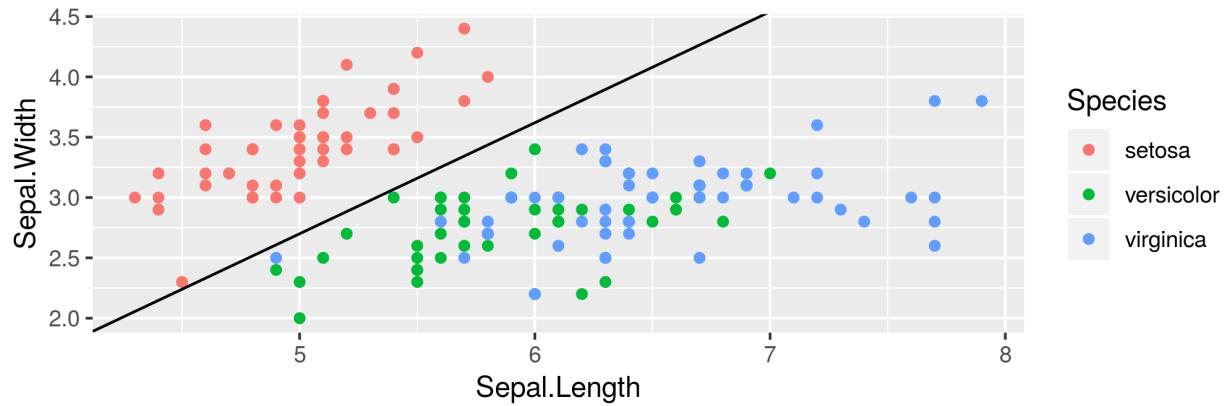
$$\Delta w_i = \eta(score_j - output_j)x_i$$

Em que x_{i_j} é o valor do i -ésimo pixel, w_i é o i -ésimo peso e η uma constante chamada *tasa de aprendizagem* (learning rate), que determina o tamanho dos incrementos feitos pelo algoritmo. Mostraremos a derivação dessa equação a seguir.

Se os dados são linearmente separáveis, o algoritmo converge com um número suficiente de exemplos.

Assim, funciona para separar flores *setosa* de outra classe, mas não teríamos bons resultados separando *virginica* de *versicolor*.

```
>ggplot(iris,aes(x=Sepal.Length,y=Sepal.Width,color=Species))+  
  geom_point() + geom_abline(slope = 0.92,intercept = -1.9)
```



Auto MaRk I

Usando as abstrações acima, codificamos nosso perceptron em R, o Auto MaRk I.

Argumentos: Exemplos (x, vetor de números reais) e estados esperados (y, disparar = 1 vs. não disparar = -1) devem ter mesmo tamanho.

Eta: Número especificando constante de aprendizagem.

Auto MaRk I inicializa um peso aleatório para cada entrada e, numa ordem aleatória, percorre os exemplos atualizando os pesos.

```
library(magrittr)
>mark_i <- function(x, y, eta) {
  # inicializa pesos randomicos de distribuicao normal
  w <- rnorm(dim(x)[2]) # numero de pesos = numero de colunas em x
  ypreds <- rep(0, dim(x)[1]) # inicializa predicoes em 0
  # Processa as observacoes em x de forma aleatoria
  for (i in sample(1:length(y), replace=F)) {
    # predicao
    ypred <- sum(w * as.numeric(x[i, ])) %>% phi_heavi
    # update em w
    delta_w <- eta * (y[i] - ypred) * as.numeric(x[i, ])
    #nota: x[i,] sera multiplicado como matriz (dot product)
    w <- w + delta_w
    ypreds[i] <- ypred # salva predicao atual
  }
  print(paste("Weights: ",w))
  return(ypreds)
}
```

Vamos testá-lo para o problema proposto, separando flores *setosa* de *versicolor*. Preparação de dados:

```
>train_df <- iris[1:100, c(1, 2, 5)]
>train_df[, 4] <- -1
>train_df[train_df[, 3] == "setosa", 4] <- 1
>names(train_df) <- c("s.len", "s.wid", "species", "target")
>head(train_df)
  s.len s.wid species target
1  5.1   3.5   setosa     1
2  4.9   3.0   setosa     1
3  4.7   3.2   setosa     1
4  4.6   3.1   setosa     1
5  5.0   3.6   setosa     1
6  5.4   3.9   setosa     1
> train_df[60:65,]
  s.len s.wid   species target
60  5.2   2.7 versicolor    -1
61  5.0   2.0 versicolor    -1
62  5.9   3.0 versicolor    -1
63  6.0   2.2 versicolor    -1
64  6.1   2.9 versicolor    -1
65  5.6   2.9 versicolor    -1
>x_features <- train_df[, c(1, 2)]
>y_target <- train_df[, 4]
```

E então, podemos ativá-lo:

Usando $\eta = 0.002$, obtivemos 41% de acurácia (classificações corretas). Podemos modificar a taxa de aprendizagem. Com $\eta = 0.05$, aumentamos para 59%. Com $\eta = 0.1$, temos 62%. Um bom valor é 0.01, com 77%. **Nada mau!** Codificamos nosso neurônio *do zero*, usando algumas matrizes, pesos aleatórios e um algoritmo sequencial de operações e atualização de pesos. Com isso, atingimos uma acurácia razoável.

```
> y_preds <- mark_i(x_features, y_target, 0.05)
[1] "Weights: -1.12748454064396"
[2] "Weights: 1.35197455996465"
> table(y_preds,train_df$target)
y_preds -1 1
-1 30 21
1 20 29

> y_preds <- mark_i(x_features, y_target, 0.1)
[1] "Weights: -2.08944843785222"
[2] "Weights: 3.35800090343738"
> table(y_preds,train_df$target)
y_preds -1 1
-1 36 14
1 14 36

> y_preds <- mark_i(x_features, y_target, 0.01)
[1] "Weights: -0.250410476080629"
[2] "Weights: 0.447470183281492"
> table(y_preds,train_df$target)
y_preds -1 1
-1 43 16
1 7 34
```

Chamamos η de hiperparâmetro. A escolha de valores para hiperparâmetros é um dos desafios em aprendizagem estatística. Uma maneira trivial é testar muitos valores possíveis e observar o desempenho, entretanto isso não é exequível para grandes volumes de dados e/ou muitos parâmetros. Existem diversos processos heurísticos e algoritmos para encontrar valores ótimos.

Uma forma popular para otimizar o treinamento é partitionar o dataset em pedaços e apresentar os particionamentos (epochs) repetidas vezes ao classificador ou acumular os erros de epochs ao invés de exemplos individuais. Assim, calculamos erros agregados e evitamos mínimos locais.

Para evitar andar em círculos, avancamos por mais tempo em uma direção antes de recalcular a rota.



Deep learning



Intuições

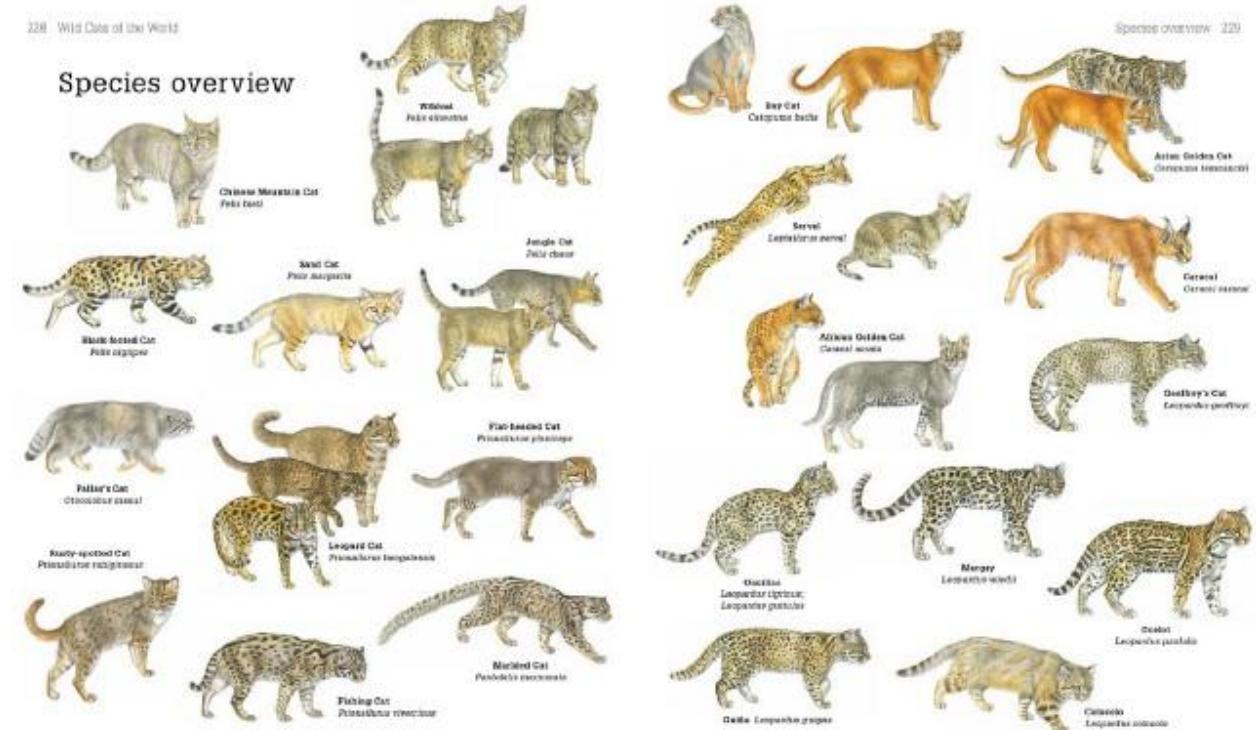
Com o aprendizado através de exemplos, otimizamos nosso classificador (mudando pesos W) para minimizar a perda usando aproximações. Assim como estendemos modelos simples do capítulo 2 usando grafos no capítulo 3, aqui vamos aplicar o mesmo conceito e imaginar relações entre neurônios.

Os dados são apresentados aos perceptrons na linha de frente. O output, porém, não é o resultado dessa primeira operação: esse valor é usado como input para neurônios da próxima camada.

Assim, conseguimos implementar transformações adequadas (rotações, torções, escalonamentos, dobradas) em sequência, de maneira que abstrações complexas possam ser capturadas.

Going Deep

As versões reais da maioria dos conceitos criados por seres humanos não são idênticas umas às outras. Em outras palavras, não existe um conjunto rígido de regras para classificarmos a maior parte das entidades ao nosso redor. Muitas entidades são diferentes, porém similares o suficiente para pertencer a uma mesma categoria.



Todos são naturalmente reconhecidos como felinos, mas apresentam variações de tamanho, cor e proporção em todo o corpo. Esse é um problema interessante e antigo, mais conhecido na ideia de entes platônicos, os quais capturam a essência de um conceito.

Alguns filósofos contemporâneos acreditam que abstrações humanas são instâncias de um conceito mais genérico: mapas biológicos contidos em redes neuronais (Paul Churchland, Plato's Camera).

Esses mapas estão associados de forma hierarquizada. Numerosos padrões em níveis inferiores e um número menor em camadas superiores.

No caso da visão, neurônios superficiais captam pontos luminosos. O padrão de ativação sensorial enviado ao córtex visual primário é o primeiro mapa, que é torcido e filtrado caminho cima.

Neurônios intermediários possuem configurações que identificam características simples: olhos e subcomponentes da face. Por fim, temos camadas mais profundas, ligadas a abstrações.

Deduzindo superfícies

Um classificador deve capturar essa estrutura abstrata a partir de modelos matemáticos tratáveis. Para examinarmos esse aspecto, usemos um exemplo. O gráfico abaixo representa milhares de amostras com: (1) a curva diária natural de um hormônio (em vermelho) e a curva sob uso de esteroides anabolizantes (azul).

Como hipotéticos membros de uma comitê atlético, nosso objetivo aqui é, dada uma amostra, saber se o atleta está sob efeito de esteroides. Quando experimentamos, normalmente haverá ruídos (erros) na medida e receberemos medições imprecisas da curva. Variações na dieta daquele dia, micções, sudorese, stress e outros fatores.

Usamos o tempo (t , eixo horizontal) e nível hormonal (β , eixo vertical).

Um modelo bastante popular para classificações é o de regressão logística. Nele, estimamos probabilidade para um evento com base nas probabilidades de uma função sigmoide. Temos uma probabilidade (valor entre 0 e 1) definida por:

$$P(h, \beta) = \frac{1}{1 + e^{-(i+t*h+\beta*y+\epsilon)}}$$

ϵ representa o erro e i é uma constante.

Em uma linha de R:

```
>logist.fit <- glm(type_dic ~ beta + tempo, family=binomial, data=inv.ds)
```

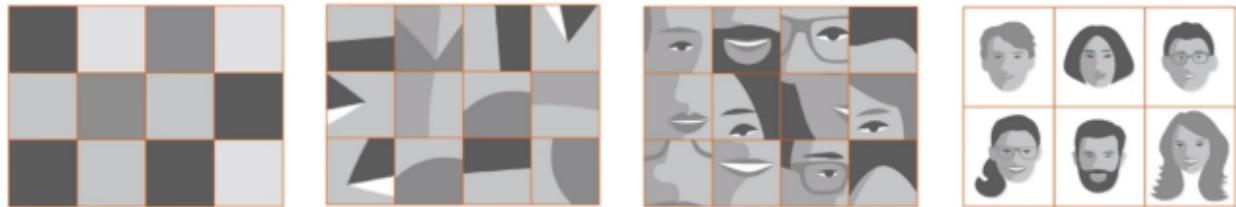
A vantagem de usar essa modelagem é que temos uma relação direta entre o inverso dessa função (P^{-1}), “logito”) e a combinação linear dos nossos parâmetros:

$$\text{logit}(P(x)) = i + t * x + \beta * y + \epsilon$$

Em outras palavras, o processo de estimação é parecido com o da regressão linear, que é facilmente tratável. Outra consequência é que assumimos que a distinção entre classes (com base no logito, log odds) pode ser dada por um limite. Este tem uma relação linear com nossas variáveis. Estimamos a magnitude e o sentido dessas relações pelos parâmetros da regressão.



Figure 6: Resposta a estímulos visuais em V1 de *Macaca fascicularis* <http://www.jneurosci.org/content/32/40/13971>



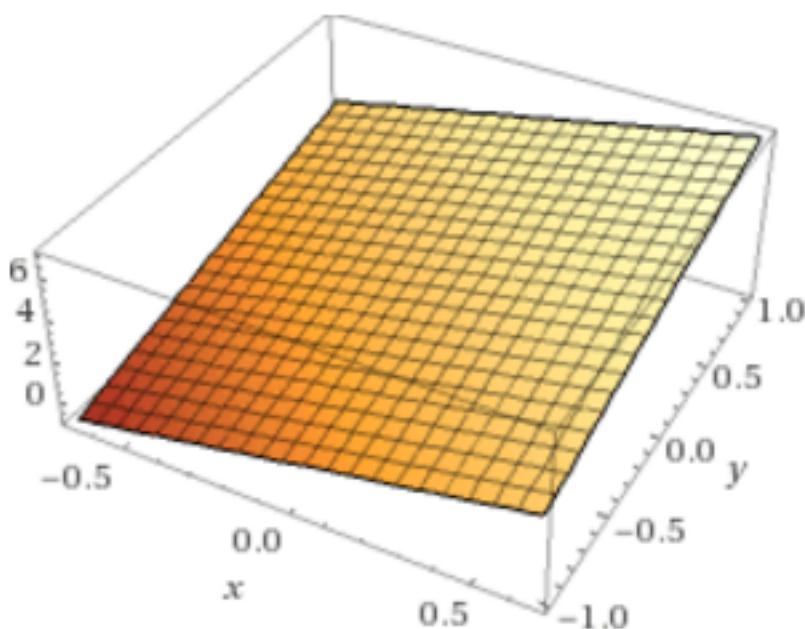
Pixels

Edges and Contrasts

Parts

Full Faces

Figure 7: Retirado de: <https://www.youtube.com/watch?v=SeyIg6ArS4Y>



Podemos imaginar que o log odds (z , eixo vertical) cresce linearmente com uma combinação de duas variáveis (x e y). Notem que a superfície definida pelo nossa equação/modelo é um [hiper]plano³ dado por $z = 3 + 3x + 2y$. Estimamos qual seria a posição na reta dada por aquela medida e usamos um limite de decisão (decision boundary) linear. Voltando ao nosso exemplo, seria difícil capturar as diferenças usando esta estratégia.

Acima, um neurônio sigmoide, que equivale à regressão logística. É como o plano anterior, mas visto de cima, dividimos ele em duas regiões para classificação. Por que? O classificador linear otimiza suas respostas levando em conta apenas o valor absoluto da medida hormonal. Isto é, valores acima de um limite serão considerados doping, não considerando horário. Matematicamente, o coeficiente para o tempo foi ajustado em 0. Mudar isso tornaria a reta divisória inclinada em relação ao eixo x , piorando a classificação.

Podemos verificar isso diretamente através dos parâmetros estimados em nosso modelo de regressão.

```
> summary(logist.fit)
Call: glm(formula = type_dic ~ beta + tempo, family = binomial, data = inv.ds)
Coefficients:
(Intercept)      beta      tempo
```

³Um hiperplano é a generalização de plano (curvatura zero) para quaisquer dimensões. O hiperplano é um espaço de $n - 1$ em um espaço n dimensões. A reta é um hiperplano em duas dimensões, o plano é um hiperplano em 3 dimensões. Para dimensões mais altas, a visualização é menos simples. Um hiperplano divide os espaço em duas partes.

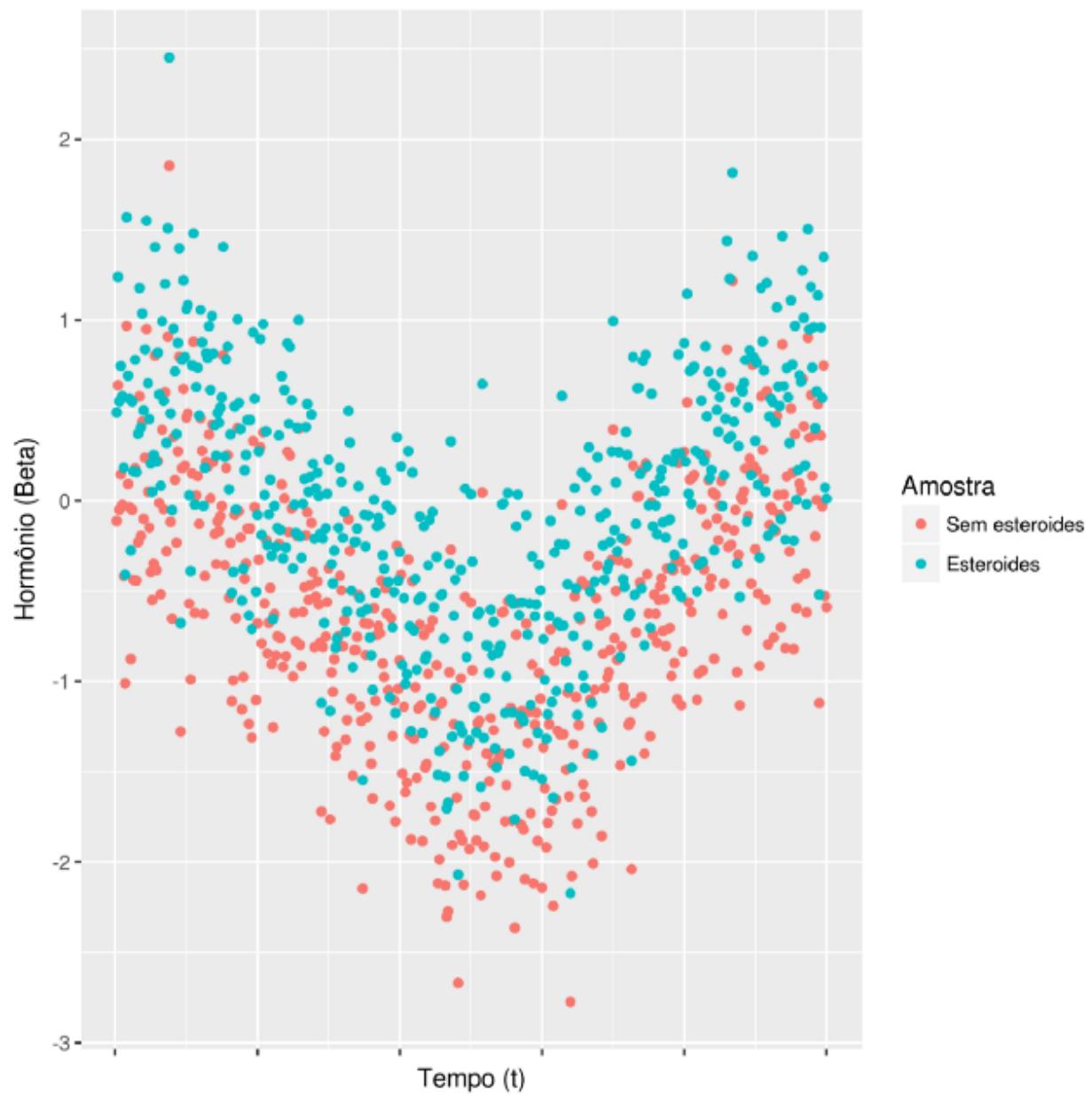


Figure 8: Exemplo inspirado no texto de Chris Olah (<http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/>)

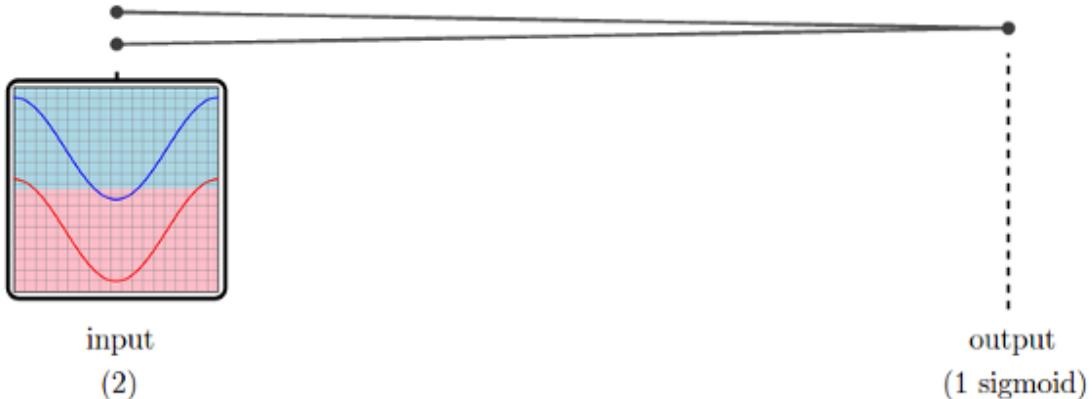


Figure 9: <http://colah.github.io/posts/2015-01-Visualizing-Representations/>

```

-0.8752803  -3.6195723  -0.0001221 # Próximo a zero
Degrees of Freedom: 999 Total (i.e. Null);  997 Residual
Null Deviance:      1386
Residual Deviance: 774.4  AIC: 780.4
> prob=predict(logist.fit,type=c("response"))
> inv.ds$prob=prob
> curve <- roc(type_dic ~ prob, data = inv.ds)
> curve

Call:
roc.formula(formula = type_dic ~ prob, data = inv.ds)
Data: prob in 500 controls (type_dic 0) < 500 cases (type_dic 1).
Area under the curve: 0.8767

```

Quem poderá nos ajudar?

A solução é introduzir termos polinomiais de grau mais alto ($x^2, x^3\dots$), interações ou usar funções mais complexas. Aí corremos o risco de realizar sobre ajuste. Deixar o sinal dos confundir e fazer um modelo complexo que não funciona em novos exemplos.

E o que acontece se conectarmos classificadores simples hierarquicamente?

A resposta de uma unidade é usada como a entrada de outras. Quando processamos o sinal em etapas, cada camada modifica os dados para as camadas posteriores, transformando e filtrando/dando forma.

As camadas intermediárias permitem a transformação gradual do sinal, e o sistema acerta usando apenas dois classificadores simples (sigmoids). No exemplo acima, temos uma camada de 2 neurônios entre o input e o output.

Agora, a primeira camada (hidden) modifica a entrada com duas unidades sigmoids e a segunda camada pode classificar corretamente usando apenas uma reta, algo que era impossível antes. Em tese, esse modelo pode capturar melhor as características que geraram os dados (flutuação hormonal ao longo do dia).

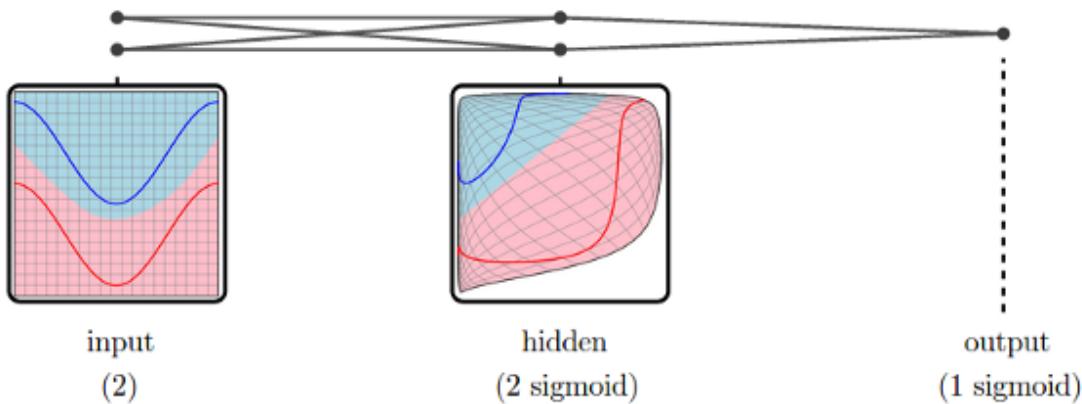


Figure 10: <http://colah.github.io/posts/2015-01-Visualizing-Representations/>

Neurônios

Notem que o diagrama acima lembra uma rede neural. Esse tipo de classificador foi inspirado na organização microscópica de neurônios reais e acredita-se que seu funcionamento seja de alguma forma análogo. A arquitetura de redes convolucionais (convolutional neural networks), estado da arte em reconhecimento de imagens, foi inspirada no córtex visual de mamíferos⁴. Outros modelos bio inspirados (Spiking neural networks, LTSMs...) apresentam desempenhos inéditos para tarefas complexas e pouco estruturadas, como reconhecimento de voz e tradução de textos. A teoria mais aceita é de que o maquinário neural dos animais foi desenhado por processos evolutivos, como a seleção natural. Assim, apresenta coloridas formas de complexidade a depender da tarefa desempenhada.

Como podemos ver, as redes biológicas são complexas, com até dezenas de bilhões de unidades de processamento paralelas conectadas. Zona destacada possui grafo isomorfo ao descrito no texto.

Nos modelos profundos (deep) de reconhecimento de rosto, neurônios de camadas superficiais capturam bordas, ângulos e vértices, camadas intermediárias detectam presença de olhos, boca, nariz. Por fim, camadas ao final da arquitetura decidem se é um rosto ou não e a quem ele pertence.

Eficiência e aplicações

Podemos demonstrar formalmente que uma rede neural com apenas uma camada interna é capaz de aproximar qualquer função. A prova não é lá essas coisas, já que, no fundo o que fazemos é criar uma tabela de consulta (lookup table) para os valores de entrada e saída usando os neurônios. Na prática, é difícil obter boas performances. Tão difícil que as redes neurais também passaram décadas esquecidas. Se você rodar o modelo abaixo, baseado no nosso exemplo, verá que a acurácia é próxima da regressão logística. É necessário algum conhecimento e tempo para afinar os detalhes.

Normalmente, depende da qualidade e da quantidade dos dados. O boom veio com a descoberta de topologias de rede especificamente boas para certas tarefas (e.g. LSTM para linguagem natural, Conv Nets para visão computacional).

Em outras palavras, modelar uma rede neural para problemas inéditos pode ser algo desafiador.

O código a seguir mostra como implementar uma rede com a topologia descrita usando a lib **deepnet**. O modelo tem um funcionamento ligeiramente diferente (deep belief networks), porém não nos preocupemos com isso no momento.

⁴(<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC1557912/>)

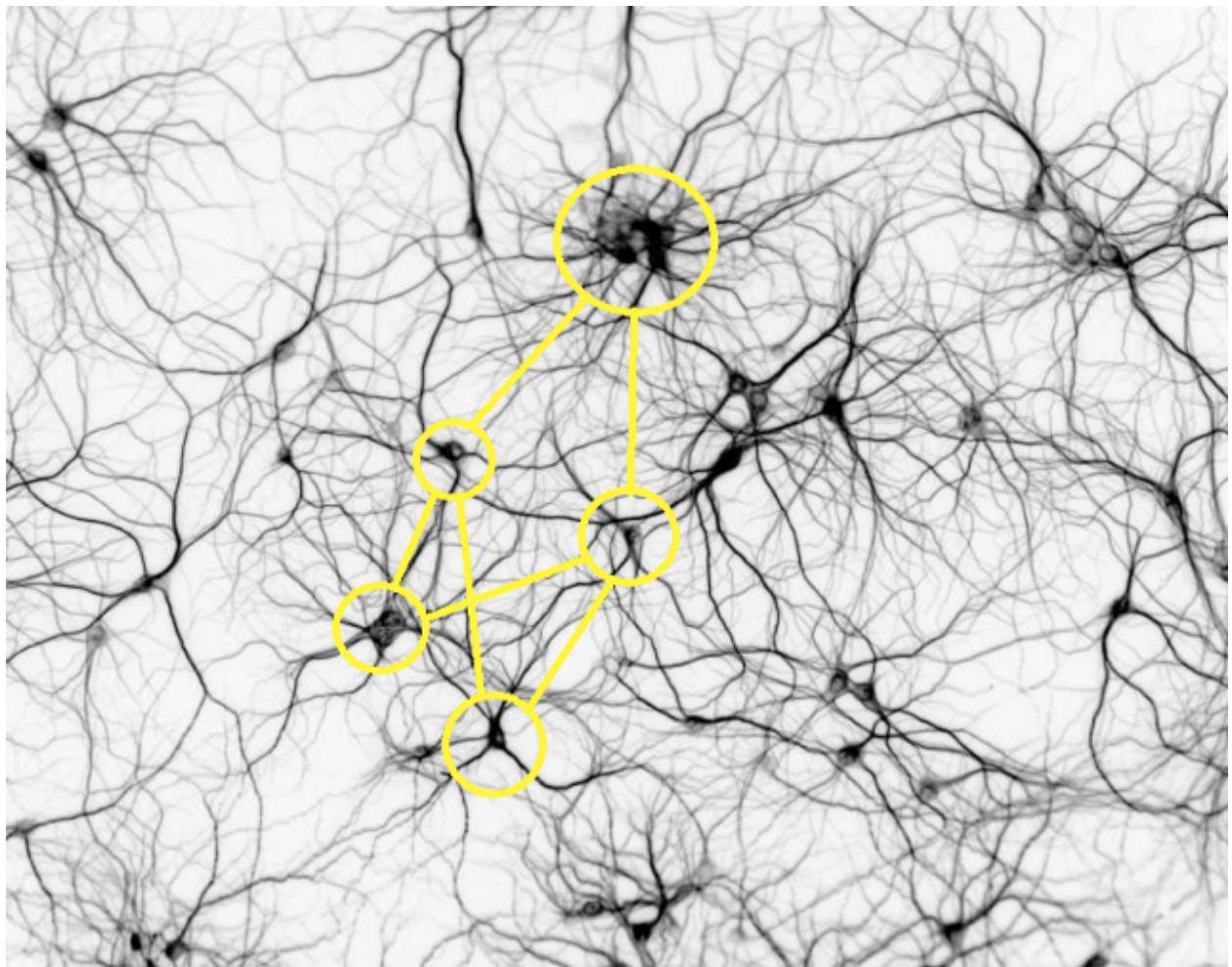


Figure 11: Modificado de <http://www.rzagabe.com/2014/11/03/an-introduction-to-artificial-neural-networks.html>

```

# Neural Net para o exemplo
# Processo para gerar dados em arquivos auxiliares ao livro(/aux)
>library(deepnet)
>inv.ds$tempo.norm <- normalize(inv.ds$tempo)
>deep.log.dbn <- dbn.dnn.train(
  x=as.matrix(inv.ds[,c("beta","tempo.norm")]),
  y=as.numeric(as.character(inv.ds$type_dic)),
  hidden = c(2), activationfun = "sigm",
  learningrate=2.65, momentum=0.85, learningrate_scale=1,
  output = "sigm", numepochs=3, batchsize= 11)
  (...)

begin to train dbn .....
training layer 1 rbm ...
dbn has been trained.

begin to train deep nn .....
deep nn has been trained.

>inv.ds$deep.test <- nn.predict(deep.log.dbn,
  x=as.matrix(inv.ds[,c("beta","tempo")]))


>curve <- roc(type_dic ~ deep.test, data=inv.ds)
>plot(curve)
>curve
Call:
roc.formula(formula = type_dic ~ deep.test, data = inv.ds)
Data: deep.test in 1000 controls (type_dic 0) < 1000 cases (type_dic 1).
Area under the curve: 0.6671

```

As redes neurais passaram algum tempo esquecidas, até que algumas reviravoltas⁵ permitiram o treinamento eficaz dessas redes. Algoritmos para melhorar o treinamento, assim como arquiteturas econômicas ou especialmente boas em determinadas tarefas. Além disso, o uso de processadores gráficos (GPU), desenhados para as operações de álgebra linear que discutimos (com matrizes) permitiu treinar em um volume maior de dados.

Gradient Descent para o Perceptron

Até o momento, ilustramos intuições e aplicações básicas, porém o grande desafio de modelos com muitos parâmetros está em orquestrar o treinamento conjunto de diferentes nodos.

Usamos os pesos com uma fórmula contendo taxa de aprendizagem (η) e outros parâmetros: a função de erro entre score desejado($score_j$) e output $E = d(score_j, output_j)$; valor da entrada (x_i).

$$w'_i = w_i + \Delta w_i$$

Δw_i pode ser obtido usando o conceito de Gradient Descent. Intuitivamente, calculamos a inclinação local e caminhamos no sentido oposto ao mais íngreme. O valor de η governa o tamanho dos passos.

Levando em conta cada j -ésima observação, definimos uma função de perda L expressando a soma dos erros nos n exemplos e minimizamos ela.

$$\min(L) = \min \sum_j^n E(score_j, output_j))$$

⁵(<http://people.idsia.ch/~juergen/who-invented-backpropagation.html>)

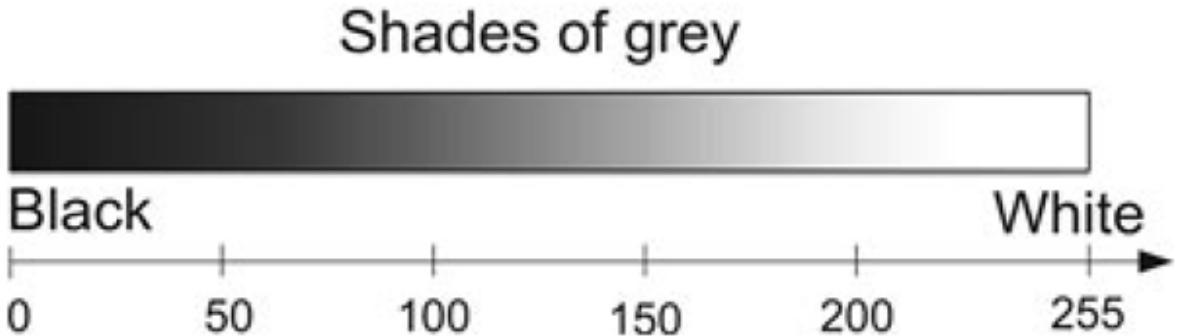
Calculamos o valor dos pesos atuais e percorremos o espaço em direção a um valor mínimo local. Se a superfície L for convexa, acharemos uma solução ótima com o número suficiente de passos.

Usaremos para nossa função de erro a distância euclidiana entre score desejado e output. O score desejado é a resposta ótima e o output é um produto entre pesos e entrada:

$$E = d_{eucl.}(score_j, output_j) = (score_j - output_j)^2$$

Notem que o processo envolve implementar uma função de erro entre resultados da rede e um espaço virtual de scores ótimos. O sucesso do treinamento depende de uma correspondência entre a função de distância escolhida e a distância real no espaço em que os dados foram gerados. Não sabemos se isso reflete a realidade. No exemplo, cada pixel reflete um sinal de 0 a 255.

A figura abaixo mostra a correspondência entre valores da medida e escala visual.



A intuição para sensibilidade à luz pode ser percebida num intervalo contínuo entre incidência total de luz (valores extremos de branco, medida: 255) e ausência total (valores extremos de preto, medida: 0). Supondo que podemos atribuir um rótulo a cada tom de cinza e que esse conjunto é ordenável pela *clareza*, dizemos que há isomorfismo de ordem entre os conjuntos.

Isso implica que a distância euclidiana deve funcionar em nossas medidas como nos números reais \mathbf{R} . Resta saber se a projeção das observações é linearmente separável. É intuitivo para seres humanos saber quais problemas serão separáveis: basta imaginar a tarefa de diferenciar tipos de imagens com uma regua numa tela em preto e branco.

Para descobrir o valor mínimo de L , vamos encontrar polos através de derivadas parciais. Ou, seu equivalente para funções de múltiplas variáveis (espaços multidimensionais), o gradiente(∇). É o produto escalar das derivadas parciais daquela função.

Para cada observação x_j , a derivada parcial da função de perda em relação a um peso w_i expressa a taxa de variação no erro global em função daquele peso. $\frac{d}{dw_i} L(w_i) = \frac{d}{dw_i} \frac{1}{n} \sum_j nE(score_j, output_j)$

Sabemos então se devemos ajustar o peso para cima ou para baixo, assim com a magnitude do passo. Algebricamente, modificaremos w seguindo o inverso do gradiente. A taxa de aprendizagem é um hiperparâmetro que regula artificialmente o tamanho desse passo:

$$\begin{aligned} \Delta w_i &= -\eta \frac{dL}{dw_i} \\ &= -\eta \frac{d}{dw_i} \frac{1}{n} \sum_j E(score_j, output_j) \end{aligned}$$

Lembrando que o erro é dado pela distância euclidiana:

$$= -\eta \frac{d}{dw_i} \frac{1}{n} \sum_j^n (score_j - output_j)^2$$

Fazemos $f(x) = (score_j - output_j)$ e $g(x) = x^2$, de maneira que

$$L = \frac{1}{n} \sum_j^n E(score_j, output_j) = (g \circ f)$$

$$= \frac{1}{n} \sum_j^n (score_j - output_j)^2$$

Podemos resolver $\frac{d}{dw_i} L$ aplicando a regra de cadeia $(g \circ f)' = (g' \circ f)f'$ e a ‘regra do tombo’ para derivadas de polinômios ($\frac{d}{dx}(x^n) = nx^{n-1}$).

Então,

$$f' = \frac{d}{dw_i} (score_j - output_j)$$

O output é dado pelo produto escalar entre pesos w_j e entradas x_j :

$$f' = \frac{d}{dw_i} (score_j - w_j \cdot x_j)$$

O score desejado não depende dos pesos, portanto a primeira derivativa é 0.

$$\begin{aligned} f' &= 0 - \frac{d}{dw_i} w_j \cdot x_j \\ &= -\frac{d}{dw_i} \sum_{i,j}^n w_{i,j} * x_{i,j} \\ &= -\frac{d}{dw_i} (w_0 * x_0 + \dots + w_i * x_i + w_n * x_n) \end{aligned}$$

Os termos não dependentes de w_i também são zerados e ficamos com o primeiro termo da soma:

$$f' = -\frac{d}{dw_i} w_i x_i$$

A função a ser derivada agora descreve uma relação linear (polinômio de grau 1) em w_i e temos:

$$f' = (-x_{i,j})$$

Sabendo f' , buscamos o outro termo em $(g \circ f)'$:

$$(g \circ f) = (score_j - output_j)^{2-1}$$

$$(g' \circ f) = 2(score_j - output_j)^{2-1}$$

$$= 2(score_j - output_j)$$

Por fim, a derivada parcial da função de perda para o i-ésimo peso w_i é:

$$\frac{dL}{dw_i} = \sum_j^n \frac{d}{dw_i} (score_j - output_j)^2$$

$$= \sum_{i,j}^n 2(score_j - w_j \cdot x_j)(-x_{i,j})$$

Para simplificar a expressão e estabelecer o tamanho dos incrementos sobre o pesos, escalamos a derivada parcial por uma constante, dada por $-\frac{1}{2}\eta_0$:

$$-\frac{1}{2} * \eta_0 \frac{dL}{dw_i} = -\frac{1}{2} \eta_0 * 2(score_j - output_j)(-x_j)$$

$$\Delta w_i = \eta_0 \sum_j^n (score_j - w \cdot x)(x_j)$$

E η_0 é um [hiper]parâmetro que simplifico a equação e define o tamanho dos incrementos usados.

Como implementamos antes no Auto MaRK I.

```
(...)
ypred <- sum(w * as.numeric(x[i, ])) %>% phi_heavi
delta_w <- eta * (y[i] - ypred) * as.numeric(x[i, ]) #-----
w <- w + delta_w
(...)
```

Backpropagation

Uma vez que o texto é sobre deep learning, precisamos falar de backpropagation. É o conceito de propagar gradientes da função perda ao longo da rede de maneira a atualizar cada nodo de maneira única.

Como vimos, podemos encarar a rede neural como uma sequência de funções plugadas. Algebraicamente, se o primeiro nodo é $q(x, y)$, o neurônio f que recebe sua saída como input tem valor $f(q(x, y))$ ou $f \circ q$.

Exemplo

Neurônio de input: $q(x, y) = 3x + 2y$

Segundo neurônio: $f(z) = z^2$

Output final: $f(q(x, y)) = q^2 = (3x + 2y)^2$

Podemos calcular o efeito de mudanças inter nodos com a regra de cadeia funções compostas. Isto é, podemos obter o gradiente de erro no nodo de hierarquia mais alta (f), com respeito a uma das variáveis de entrada (x) na hierarquia mais baixa. A operação é computacionalmente barata, bastando multiplicar as derivadas parciais dos erros em cada parte.

$$\frac{df}{dx} = \frac{df}{dq} \frac{dq}{dx}$$

É possível calcular de forma recursiva, portanto local e paralela, ao longo das camadas. Fazendo o mesmo acima para df/dy , teremos os valores de df/dx e df/dy que é precisamente nosso gradiente em f .

```
# Valor duplo (x,y) para inputs
>x <- 1
>y <- 3
q <- 3*x + 2*y # primeira camada
f <- q^2 # segunda camada
```

```

# Backprop - Mudanças em hierarquia superior
# dadas por entradas de camadas inferiores
dfdq <- 2*q # derivada de  $x^2$ ; variação de  $f$  em função de  $q$ 
dqdx <- 3 # Derivada de  $3x$ ; variação de  $q$  em função de  $x$ 
dqdy <- 2 # Derivada de  $2x$ ; variação de  $q$  em função de  $y$ 
# Obter gradiente de  $f(x,y)$  multiplicando as parciais
dfdx = dfdq*dqdx
dfdy = dfdq*dqdy
grad = c(dfdx,dfdy)
> grad
[1] 24 16

```

Usando essa lógica, calculamos os gradientes para a função de erro e treinamos o modelo.

Podemos então implementar nossa rede neural, Mark II.

Mark II

Nossa rede terá um perceptron de entrada com dimensão igual à do input. Entretanto, acrescentamos um peso a mais, que corresponderá a um intercepto.

Note que

$$y = w_0 + w_1x_1 + w_2x_2$$

é o mesmo que

$$y = 1 * w_0 + w_1x_1 + w_2x_2$$

Assim, adicionamos um peso w_0 e também forçamos uma dimensão a mais no input, que sempre terá valor 1. Chamamos esse artifício de adicionar um intercepto (*bias*) de *bias trick*. Ajuda a estabelecer um valor basal para o output, facilitando a convergência.

```

library(magrittr)
library(ggplot2)
set.seed(2600)

mark_ii <- function(x, y, eta, reps=1) {

  # inicializa pesos randomicos de distribuicao normal
  w1 <- rnorm(n = (dim(x)[2]+1)) %>% as.matrix # numero de pesos = numero de colunas em x + bias

```

Em seguida, neurônios da camada intermediária, dois, cada um com dois pesos.

```

w21 <- rnorm(2) %>% as.matrix
w22 <- rnorm(2) %>% as.matrix

```

Zeramos as previsões e iniciamos os loops de treinamento. Para a rede neural, precisamos de muitos exemplos de exposição, então embutimos em Mark II um parâmetro (*reps*) responsável por repetir o processo de treinamento com o dataset.

A rigor o melhor seria partitionar o dataset em fragmentos menores para cada epoch, mas vamos manter as coisas simples.

```

ypreds <- rep(0,dim(x)[1]) # inicializa predicoes em 0
yerrors <- rep(0,dim(x)[1]) # inicializa predicoes em 0
for (j in 1:reps){
  print(paste("This is training epoch:",j))
  print(paste("Current weights:",w1,w21,w21))

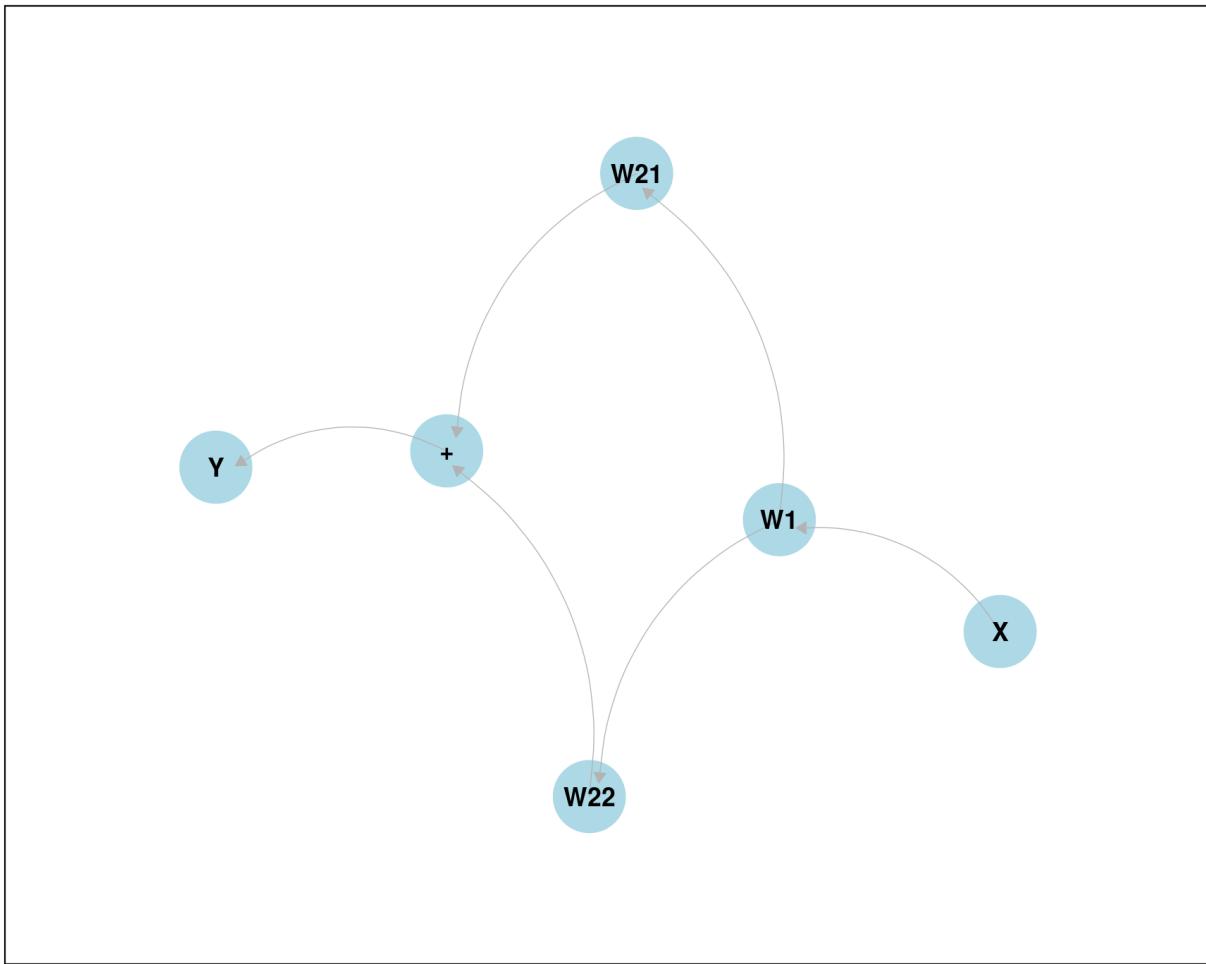
```

Predições: a primeira camada soma o produto de seus pesos pela entrada e pela unidade (*bias trick*). Os neurônios da segunda camada somam o produto de seus pesos pelo output. O output final é a soma dos outputs na camada intermediária.

```
# Processa as observações em x de forma aleatoria
for (i in sample(1:length(y), replace=F)) {
  # predicao
  ypred1 <- sum(w1 %*% c(as.numeric(x[i, ]), 1))

  ypred21 <- sum(w21 %*% as.numeric(ypred1))
  ypred22 <- sum(w22 %*% as.numeric(ypred1))

  out <- sum(ypred21, ypred22)
```



Agora, as regras de atualização dos pesos seguindo derivações com regra de cadeia. Para os neurônios intermediários, temos: $\frac{d}{dw_{21}}$ e $\frac{d}{dw_{21}}$ de $(target - (pred22 + pred21))^2$.

$$\frac{d}{dw_{21}} (target - (pred22 + pred21))^2$$

Aplicando a regra de cadeia e sabendo que a predição do segundo neurônio W_{22} não depende dos pesos em W_{21} :

$$\begin{aligned}
&= 2(\text{target} - (\text{pred22} + \text{pred21})) * \frac{d}{dw_{21}}(-1)(\text{pred22} + \text{pred21}) \\
&= 2(\text{target} - (\text{pred22} + \text{pred21})) * \frac{d}{dw_{21}}(-1)(w_{21} * \text{ypred1})
\end{aligned}$$

Que é a derivada para os pesos do perceptron:

$$= 2(\text{target} - (\text{pred22} + \text{pred21})) * (\text{ypred1})(-1)$$

Entretanto, calcular os pesos de w_1 em função da saída requer um pouco mais:

$$\begin{aligned}
&\frac{d}{dw_1}(\text{target} - (\text{pred22} + \text{pred21}))^2 \\
&= 2(\text{target} - (\text{pred22} + \text{pred21})) * \frac{d}{dw_1}(-1)(\text{pred22} + \text{pred21}) \\
&= 2(\text{target} - (\text{pred22} + \text{pred21})) * \frac{d}{dw_1}(-1)(\sum w_{22} \sum w_1 x + \sum w_{21} \sum w_1 x)
\end{aligned}$$

Usando a derivada de somas e verificando que os termos não relacionados ao w_1 avaliado somem:

$$2(\text{target} - (\text{pred22} + \text{pred21})) * (-1)(\sum w_{22}x + \sum w_{21}x)$$

```

# update em w . Eta ja ajustado para 1/2*eta
delta_w22 <- eta * (-1) * (y[i] - (ypred21 + ypred22)) * ypred1
delta_w21 <- eta * (-1) * (y[i] - (ypred21 + ypred22)) * ypred1
delta_w1 <- eta * (y[i] - (ypred21 + ypred22)) * -1 *
  (sum(w21 %*% c(as.numeric(x[i,]),1)) + sum(w22 %*% c(as.numeric(x[i,]),1)))

w1 <- w1 - delta_w1
w21 <- w21 - delta_w21
w22 <- w22 - delta_w22
ypreds[i] <- out # salva predicao21 atual
yerrors[i] <- ypreds[i] - y[i]
}
print(paste("Mean squared error:", mean((yerrors)^2)))
}
return(ypreds)
}

```

Então, podemos testá-lo em um dataset.

```

>train_df <- iris[, c(1, 2, 3)]
>names(train_df) <- c("s.len", "s.wid", "p.len")
>head(train_df)
>train_df[60:65,]

>x_features <- train_df[, c(1, 2)]
>y_target <- train_df[, 3]

# Convergência boa
>mark_ii_preds <- mark_ii(x = x_features, y = y_target,
                           eta=0.000001, reps = 40)

```

```

[1] "This is training epoch: 1"
[1] "Current weights: -0.45050790019773 -0.0197893400687895 -0.0197893400687895"
[2] "Current weights: 0.150011803623929 2.13458518518008 2.13458518518008"
[3] "Current weights: 1.48235899015804 -0.0197893400687895 -0.0197893400687895"
[1] "Mean squared error: 1133.22204821886"
(...)

[1] "This is training epoch: 2"
[1] "Current weights: -0.67126807499406 -0.0609395311239563 -0.0609395311239563"
[2] "Current weights: -0.0707483711724013 2.09343499412492 2.09343499412492"
[3] "Current weights: 1.26159881536171 -0.0609395311239563 -0.0609395311239563"
[1] "Mean squared error: 176.747586724131"
(...)

[1] "This is training epoch: 4"
[1] "Current weights: -0.791488817323548 -0.0700721883119202 -0.0700721883119202"
[2] "Current weights: -0.19096911350189 2.08430233693696 2.08430233693696"
[3] "Current weights: 1.14137807303222 -0.0700721883119202 -0.0700721883119202"
[1] "Mean squared error: 7.32496712284895"
[1] "This is training epoch: 5"
[1] "Current weights: -0.805708526415977 -0.0705118739404967 -0.0705118739404967"
[2] "Current weights: -0.205188822594319 2.08386265130838 2.08386265130838"
[3] "Current weights: 1.12715836393979 -0.0705118739404967 -0.0705118739404967"
[1] "Mean squared error: 3.31246116798174"
(...)
[1] "Mean squared error: 2.50706426321967"
(...)
[1] "Mean squared error: 2.50638029884829"
(...)
[1] "Mean squared error: 2.50640582517322"

```

Podemos observar o modelo convergindo à medida em que os pesos se estabilizam e nossa medida de erro cai. Usando o η acima, a rede costuma convergir com correlação $\rho \sim 0.60$ entre previsões e dados originais.

```

>acc_data <- data.frame(y_preds=mark_ii_preds,
                           y_targs=y_target)

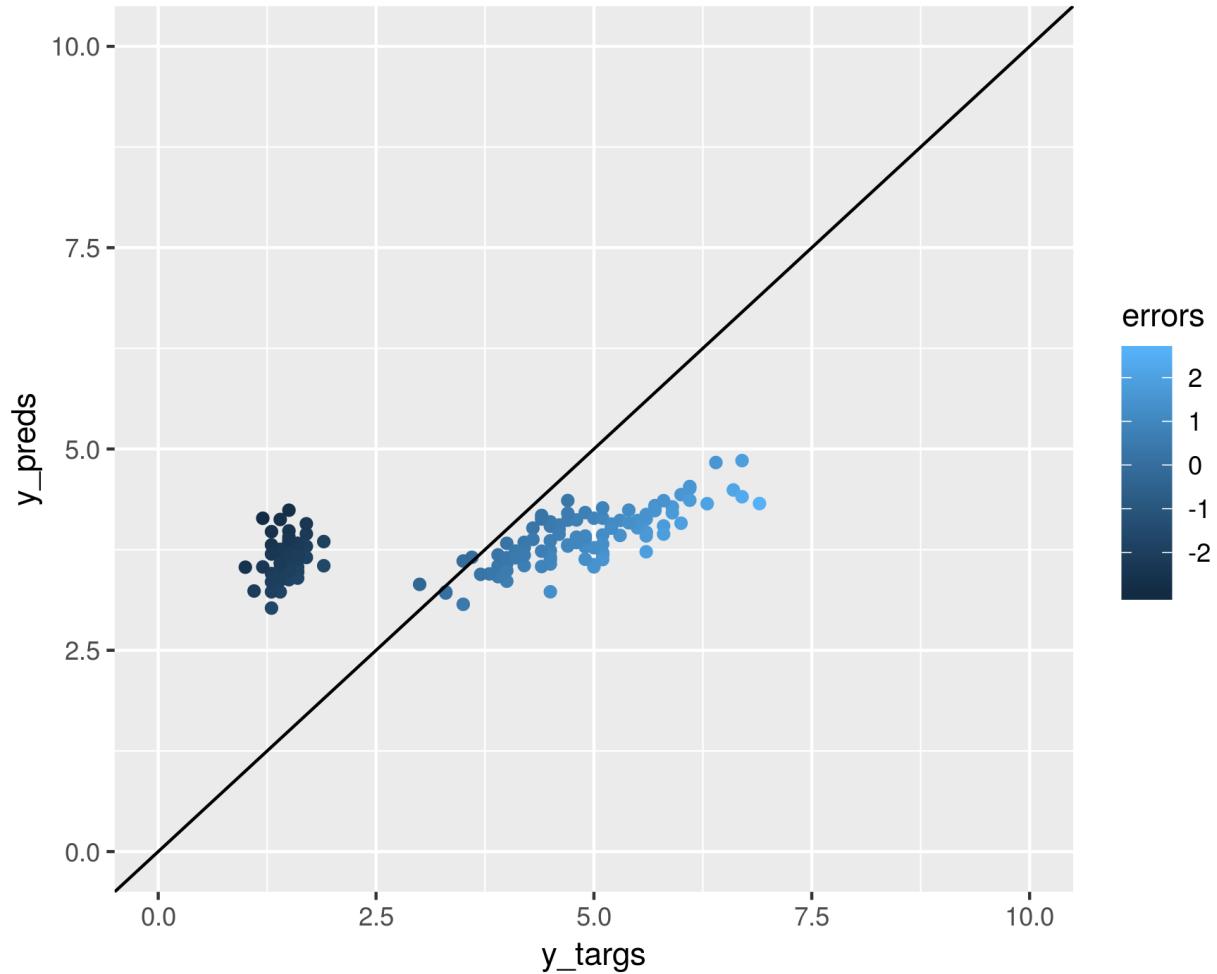
>acc_data$errors <- y_target - mark_ii_preds
>cor.test(acc_data$y_preds, acc_data$y_targs)

Pearson's product-moment correlation

data: acc_data$y_preds and acc_data$y_targs
t = 8.9717, df = 148, p-value = 1.203e-15
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
0.4788098 0.6883163
sample estimates:
cor
0.5935271

>ggplot(acc_data,aes(y=y_preds,x=y_targs,color=errors))+ 
  geom_point() + xlim(0,10) + ylim(0,10) +
  geom_abline(slope = 1,intercept = 0)

```



Referências

Para uma história completa sobre redes neurais: J. Schmidhuber. Deep Learning in Neural Networks: An Overview. *Neural Networks*, 61, p 85–117, 2015. (Based on 2014 TR with 88 pages and 888 references, with PDF & LATEX source & complete public BIBTEX file).

<http://web.csulb.edu/~cwallis/artificialn/History.htm> https://sebastianraschka.com/Articles/2015_singlelayer_neurons.html <https://rpubs.com/FaiHas/197581>