линейная регрессия по-байесовски

Сергей Николенко

СПбГУ — Санкт-Петербург 30 сентября 2023 г.

Random facts:

- 30 сентября 1452 г. в Майнце Иоганн Гутенберг напечатал свою первую Библию
- 30 сентября 1846 г. Уильям Мортон впервые вырвал зуб с использованием анестезии диэтилового эфира; Мортон перенял эту идею у Хораса Уэллса, который успешно применял веселящий газ во врачебной практике, но публичная демонстрация прошла неудачно, и Уэллс покончил жизнь самоубийством в 1848 году
- 30 сентября 1882 г. в городе Эпплтон (штат Висконсин) на реке Фокс заработала первая в мире гидроэлектростанция по системе Эдисона мощностью 12.5кВт
- 30 сентября 1929 г. прошла первая телетрансляция ВВС, а 30 сентября 1939 г. NBC впервые провела телетрансляции матча по американскому футболу
- · 30 сентября 1520 г. стал султаном Сулейман I Великолепный
- 30 сентября 1938 г. было подписано Мюнхенское соглашение между Германией,
 Великобританией, Францией и Италией о передаче Судетской области Германии
- 30 сентября 2005 г. в датской газете «Jyllands-Posten» были опубликованы двенадцать карикатур на пророка Мухаммеда

Регуляризация по-байесовски —

- А теперь давайте посмотрим на регрессию с совсем байесовской стороны.
- Напомним основу байесовского подхода:
 - 1. найти апостериорное распределение на гипотезах/параметрах:

$$p(\theta \mid D) \propto p(D|\theta)p(\theta)$$

(и/или найти максимальную апостериорную гипотезу $\arg\max_{a} p(\theta \mid D)$);

2. найти апостериорное распределение исходов дальнейших экспериментов:

$$p(x \mid D) \propto \int_{\theta \in \Theta} p(x \mid \theta) p(D|\theta) p(\theta) d\theta.$$

- В нашем рассмотрении пока не было никаких априорных распределений.
- Давайте какое-нибудь введём; например, нормальное (почему так позже):

$$p(\mathbf{w}) = N(\mathbf{w} \mid \mu_0, \Sigma_0).$$

• Рассмотрим набор данных $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ со значениями $\mathbf{t} = \{t_1, \dots, t_N\}$. В этой модели мы предполагаем, что данные независимы и одинаково распределены:

$$p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N N\left(t_n \mid \mathbf{w}^\top \phi(\mathbf{x}_n), \sigma^2\right).$$

• Тогда наша задача – посчитать

$$\begin{split} p(\mathbf{w} \mid \mathbf{t}) &\propto p(\mathbf{t} \mid \mathbf{X}, \mathbf{w}, \sigma^2) p(\mathbf{w}) \\ &= N(\mathbf{w} \mid \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \prod_{n=1}^N N\left(t_n \mid \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n), \sigma^2\right). \end{split}$$

• Давайте подсчитаем.

• Получится

$$\begin{split} p(\mathbf{w} \mid \mathbf{t}) &= N(\mathbf{w} \mid \mu_N, \Sigma_N), \\ \mu_N &= \Sigma_N \left(\Sigma_0^{-1} \mu_0 + \frac{1}{\sigma^2} \boldsymbol{\Phi}^\top \mathbf{t} \right), \\ \Sigma_N &= \left(\Sigma_0^{-1} + \frac{1}{\sigma^2} \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \right)^{-1}. \end{split}$$

• Теперь давайте подсчитаем логарифм правдоподобия.

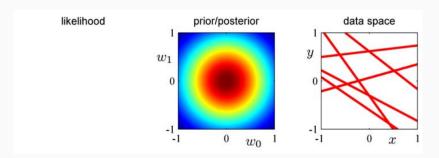
• Если мы возьмём априорное распределение около нуля:

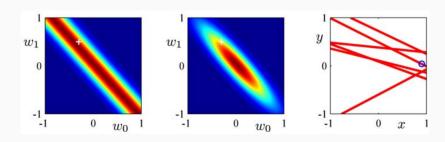
$$p(\mathbf{w}) = N(\mathbf{w} \mid 0, \frac{1}{\alpha} \mathbf{I}),$$

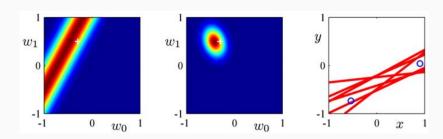
то логарифм правдоподобия получится

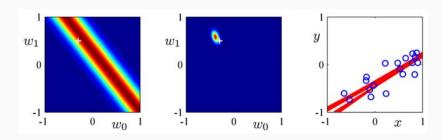
$$\ln p(\mathbf{w} \mid \mathbf{t}) = -\frac{1}{2\sigma^2} \sum_{n=1}^{N} \left(t_n - \mathbf{w}^{\top} \phi(\mathbf{x}_n) \right)^2 - \frac{\alpha}{2} \mathbf{w}^{\top} \mathbf{w} + \text{const},$$

то есть в точности гребневая регрессия.









Обобщение

 Можно слегка обобщить – рассмотреть априорное распределение более общего вида

$$p(\mathbf{w} \mid \alpha) = \left[\frac{q}{2} \left(\frac{\alpha}{2} \right)^{1/q} \frac{1}{\Gamma(1/q)} \right]^M e^{-\frac{\alpha}{2} \sum_{j=1}^{M} \left| w_j \right|^q}.$$

Упражнение. Подсчитайте логарифм правдоподобия.

Лассо

• Теперь давайте рассмотрим лассо-регрессию:

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (f(x_i, \mathbf{w}) - y_i)^2 + \lambda \sum_{j=0}^p |w_j|. \label{eq:loss}$$

- Главное отличие теперь форма ограничений (т.е. форма априорного распределения) такова, что весьма вероятно получить строго нулевые w_i .
- · Кстати, что значит «форма ограничений»?

Лассо

 Мы можем переписать регрессию с регуляризатором по-другому:

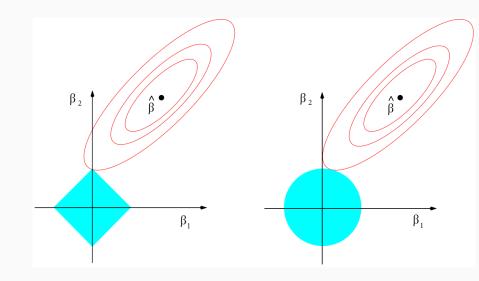
$$\mathbf{w}^* = \mathop{\arg\min}_{\mathbf{w}} \left\{ \frac{1}{2} \sum_{i=1}^N (f(x_i, \mathbf{w}) - y_i)^2 + \lambda \sum_{j=0}^p |w_j| \right\},$$

эквивалентно

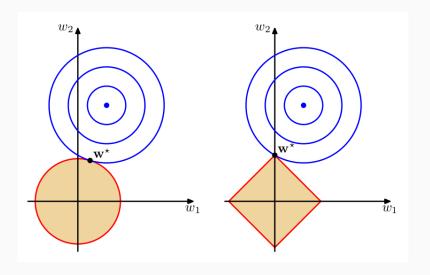
$$\mathbf{w}^* = \arg\min_{\mathbf{w}} \left\{ \frac{1}{2} \sum_{i=1}^N (f(x_i, \mathbf{w}) - y_i)^2 \right\} \text{ при } \sum_{j=0}^p |w_j| \leq t.$$

Упражнение. Докажите это.

ГРЕБЕНЬ И ЛАССО



ГРЕБЕНЬ И ЛАССО



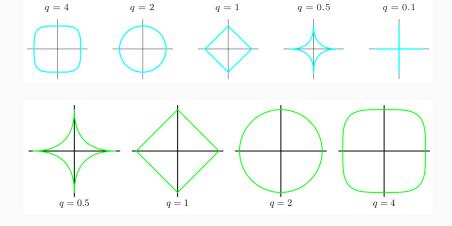
Обобщение

 Можно рассмотреть обобщение гребневой и лассо-регрессии:

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^{N} (f(x_i, \mathbf{w}) - y_i)^2 + \lambda \sum_{j=0}^{p} (|w_j|)^q.$$

Упражнение. Какому априорному распределению на параметры ${\bf w}$ соответствует эта задача?

$\mathsf{P}\mathsf{a}\mathsf{3}\mathsf{h}\mathsf{b}\mathsf{i}\mathsf{e}\ q$



В ЛИНЕЙНОЙ РЕГРЕССИИ

Предсказание в линейной регрессии

- Теперь давайте вернёмся к байесовской постановке:
 - 1. найти апостериорное распределение на гипотезах/параметрах:

$$p(\theta \mid D) \propto p(D|\theta)p(\theta)$$

(и/или найти максимальную апостериорную гипотезу $\arg\max_{\theta} p(\theta \mid D)$);

2. найти апостериорное распределение исходов дальнейших экспериментов:

$$p(x \mid D) \propto \int_{\theta \in \Theta} p(x \mid \theta) p(D|\theta) p(\theta) d\theta.$$

ПРЕДСКАЗАНИЕ В ЛИНЕЙНОЙ РЕГРЕССИИ

• В прошлый раз мы нашли апостериорное распределение: для гауссовского априорного

$$p(\mathbf{w} \mid \alpha) = N(\mathbf{w} \mid 0, \frac{1}{\alpha}\mathbf{I})$$

мы нашли

$$\begin{split} p(\mathbf{w} \mid \mathbf{t}, \alpha, \beta) &= N(\mathbf{w} \mid \mu_N, \Sigma_N), \\ \mu_N &= \Sigma_N \left(\Sigma_0^{-1} \mu_0 + \beta \boldsymbol{\Phi}^\top \mathbf{t} \right), \\ \Sigma_N &= \left(\Sigma_0^{-1} + \beta \boldsymbol{\Phi}^\top \boldsymbol{\Phi} \right)^{-1}, \end{split}$$

где $\beta=\frac{1}{\sigma^2}$ (precision нормального распределения).

Предсказание в линейной регрессии

• Теперь сделаем следующий шаг – найдём апостериорное распределение наших предсказаний

$$p(t \mid \mathbf{t}, \alpha, \beta) = \int p(t \mid \mathbf{w}, \beta) p(\mathbf{w} \mid \mathbf{t}, \alpha, \beta) d\mathbf{w}.$$

• Это свёртка двух гауссианов, и получается...

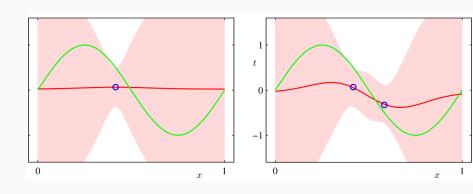
ПРЕДСКАЗАНИЕ В ЛИНЕЙНОЙ РЕГРЕССИИ

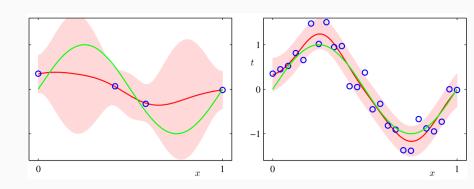
• ...тоже гауссиан:

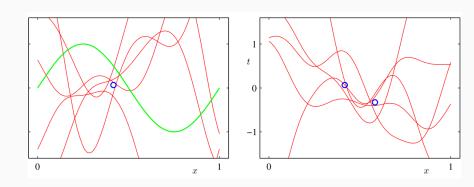
$$p(t \mid \mathbf{t}, \alpha, \beta) = N(t \mid \mu_N^{\top} \phi(\mathbf{x}), \sigma_N^2),$$
 где $\sigma_N^2 = \frac{1}{\beta} + \phi(\mathbf{x})^{\top} \Sigma_N \phi(\mathbf{x}).$

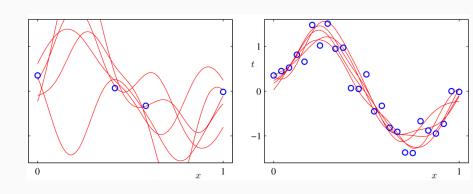
• Т.е. дисперсия складывается из шума в данных β и дисперсии параметров \mathbf{w} ; гауссианы независимы, и их дисперсии просто складываются.

Упражнение. Оценка всё время уточняется: $\sigma_{N+1}^2 \leq \sigma_N^2$.









Введение в классификацию

Задача классификации

- Теперь классификация: определить вектор ${\bf x}$ в один из K классов C_k .
- В итоге у нас так или иначе всё пространство разобьётся на эти классы.
- Т.е. на самом деле мы ищем разделяющую поверхность (decision surface, decision boundary).

Задача классификации

- Как кодировать? Бинарная задача очень естественно, переменная $t,\,t=0$ соответствует $C_1,\,t=1$ соответствует $C_2.$
- Оценку t можно интерпретировать как вероятность (по крайней мере, мы постараемся, чтобы было можно).
- Если несколько классов удобно 1-of-*K*:

$$\mathbf{t} = (0,\ldots,0,1,0,\ldots)^{\top}.$$

 Тоже можно интерпретировать как вероятности – или пропорционально им.

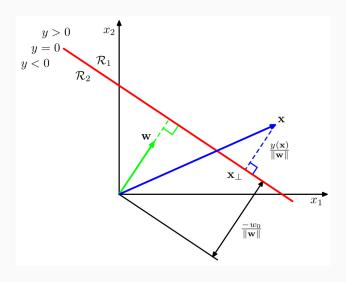
Разделяющая гиперплоскость

• Начнём с геометрии: рассмотрим линейную дискриминантную функцию

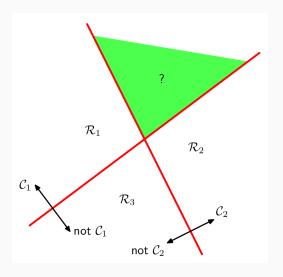
$$y(\mathbf{x}) = \mathbf{w}^{\top} \mathbf{x} + w_0.$$

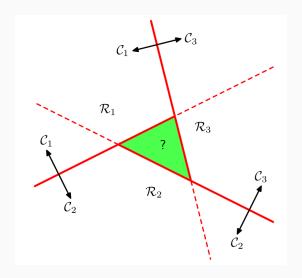
- \cdot Это гиперплоскость, и ${f w}$ нормаль к ней.
- Расстояние от начала координат до гиперплоскости равно $\frac{-w_0}{\|\mathbf{w}\|}$.
- $\cdot \ y(\mathbf{x})$ связано с расстоянием до гиперплоскости: $d = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}.$

Разделяющая гиперплоскость



- С несколькими классами выходит незадача.
- Можно рассмотреть K поверхностей вида «один против всех».
- Можно $\binom{K}{2}$ поверхностей вида «каждый против каждого».
- Но всё это как-то нехорошо.





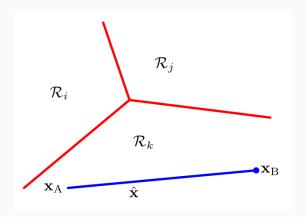
• Лучше рассмотреть единый дискриминант из K линейных функций:

$$y_k(\mathbf{x}) = \mathbf{w}_k^{\intercal} \mathbf{x} + w_{k0}.$$

- Классифицировать в C_k , если $y_k(\mathbf{x})$ максимален.
- · Тогда разделяющая поверхность между C_k и C_j будет гиперплоскостью вида $y_k(\mathbf{x})=y_j(\mathbf{x})$:

$$\left(\mathbf{w}_k - \mathbf{w}_j\right)^{\top} \mathbf{x} + \left(w_{k0} - w_{j0}\right).$$

Несколько классов



Упражнение. Докажите, что области, соответствующие классам, при таком подходе всегда односвязные и выпуклые.

МЕТОД НАИМЕНЬШИХ КВАДРАТОВ

• Мы снова можем воспользоваться методом наименьших квадратов: запишем $y_k(\mathbf{x}) = \mathbf{w}_k^{\top} \mathbf{x} + w_{k0}$ вместе (спрятав свободный член) как

$$\mathbf{y}(\mathbf{x}) = \mathbf{W}^{\top} \mathbf{x}.$$

 Можно найти W, оптимизируя сумму квадратов; функция ошибки:

$$E_D(\mathbf{W}) = \frac{1}{2} \mathrm{Tr} \left[\left(\mathbf{X} \mathbf{W} - \mathbf{T} \right)^\top \left(\mathbf{X} \mathbf{W} - \mathbf{T} \right) \right].$$

• Берём производную, решаем...

МЕТОД НАИМЕНЬШИХ КВАДРАТОВ

• ...получается привычное

$$\mathbf{W} = \left(\mathbf{X}^{\top}\mathbf{X}\right)^{-1}\mathbf{X}^{\top}\mathbf{T} = \mathbf{X}^{\dagger}\mathbf{T},$$

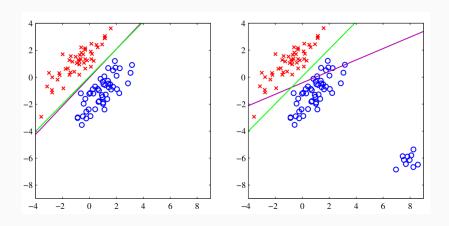
где \mathbf{X}^\dagger – псевдообратная Мура-Пенроуза.

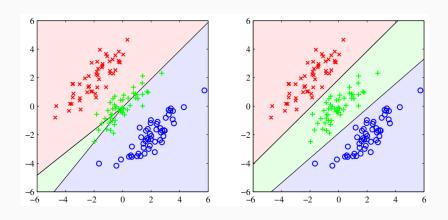
• Теперь можно найти и дискриминантную функцию:

$$\mathbf{y}(\mathbf{x}) = \mathbf{W}^{\top} \mathbf{x} = \mathbf{T}^{\top} (\mathbf{X}^{\dagger})^{\top} \mathbf{x}.$$

МЕТОД НАИМЕНЬШИХ КВАДРАТОВ

- Это решение сохраняет линейность. Упражнение. Докажите, что в схеме кодирования 1-of-K предсказания $y_k(\mathbf{x})$ для разных классов при любом \mathbf{x} будут давать в сумме 1. Почему они всё-таки не будут разумными оценками вероятностей?
 - Проблемы наименьших квадратов:
 - · outliers плохо обрабатываются;
 - · «слишком правильные» предсказания добавляют штраф.





• Почему так? Почему наименьшие квадраты так плохо работают?

- Почему так? Почему наименьшие квадраты так плохо работают?
- Они предполагают гауссовское распределение ошибки.
- Но, конечно, распределение у бинарных векторов далеко не гауссово.

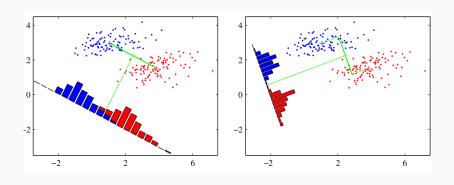
- Другой взгляд на классификацию: в линейном случае мы хотим спроецировать точки в размерность 1 (на нормаль разделяющей гиперплоскости) так, чтобы в этой размерности 1 они хорошо разделялись.
- Т.е. классификация это такой метод радикального сокращения размерности.
- Давайте посмотрим на классификацию с этих позиций и попробуем добиться оптимальности в каком-то смысле.

- Рассмотрим два класса C_1 и C_2 с N_1 и N_2 точками.
- Первая идея надо найти серединный перпендикуляр между центрами кластеров

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{C_1} \mathbf{x}, \text{ if } \mathbf{m}_2 = \frac{1}{N_2} \sum_{C_2} \mathbf{x},$$

т.е. максимизировать $\mathbf{w}^{ op}\left(\mathbf{m}_{2}-\mathbf{m}_{1}
ight)$.

• Надо ещё добавить ограничение $\|\mathbf{w}\|=1$, но всё равно не ахти как работает.



Чем левая картинка хуже правой?

- Слева больше дисперсия каждого кластера.
- Идея: минимизировать перекрытие классов, оптимизируя и проекцию расстояния, и дисперсию.
- · Выборочные дисперсии в проекции: для $y_n = \mathbf{w}^{ op} \mathbf{x}_n$

$$s_1 = \sum_{n \in C_1} \left(y_n - m_1 \right)^2 \text{ if } s_1 = \sum_{n \in C_2} \left(y_n - m_2 \right)^2.$$

• Критерий Фишера:

$$\begin{split} J(\mathbf{w}) &= \frac{\left(m_2 - m_1\right)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^\top \mathbf{S}_B \mathbf{w}}{\mathbf{w}^\top \mathbf{S}_W \mathbf{w}}, \text{ где} \\ \mathbf{S}_B &= \left(\mathbf{m}_2 - \mathbf{m}_1\right) \left(\mathbf{m}_2 - \mathbf{m}_1\right)^\top, \\ \mathbf{S}_W &= \sum_{n \in C_1} \left(\mathbf{x}_n - \mathbf{m}_1\right) \left(\mathbf{x}_n - \mathbf{m}_1\right)^\top + \sum_{n \in C_2} \left(\mathbf{x}_n - \mathbf{m}_2\right) \left(\mathbf{x}_n - \mathbf{m}_2\right)^\top. \end{split}$$

(between-class covariance и within-class covariance).

• Дифференцируя по w...

 \cdot ...получим, что $J(\mathbf{w})$ максимален при

$$\left(\mathbf{w}^{\top}\mathbf{S}_{B}\mathbf{w}\right)\mathbf{S}_{W}\mathbf{w}=\left(\mathbf{w}^{\top}\mathbf{S}_{W}\mathbf{w}\right)\mathbf{S}_{B}\mathbf{w}.$$

- \cdot Т.к. $\mathbf{S}_B = (\mathbf{m}_2 \mathbf{m}_1) \left(\mathbf{m}_2 \mathbf{m}_1\right)^{\mathsf{T}}$, $\mathbf{S}_B \mathbf{w}$ всё равно будет в направлении $\mathbf{m}_2 \mathbf{m}_1$, а длина \mathbf{w} нас не интересует.
- Поэтому получается

$$\mathbf{w} \propto \mathbf{S}_W^{-1} \left(\mathbf{m}_2 - \mathbf{m}_1 \right).$$

• В итоге мы выбрали направление проекции, и осталось только разделить данные на этой проекции.

- Любопытно, что дискриминант Фишера тоже можно получить из наименьших квадратов.
- Давайте для класса C_1 выберем целевое значение $\frac{N_1+N_2}{N_1}$, а для класса C_2 возьмём $-\frac{N_1+N_2}{N_2}$.

Упражнение. Докажите, что при таких целевых значениях наименьшие квадраты

– это дискриминант Фишера.

· А что будет с несколькими классами? Рассмотрим $\mathbf{y} = \mathbf{W}^{\top}\mathbf{x}$, обобщим внутреннюю дисперсию как

$$\mathbf{S}_W = \sum_{k=1}^K \mathbf{S}_k = \sum_{k=1}^K \sum_{n \in C_k} \left(\mathbf{x}_n - \mathbf{m}_k\right) \left(\mathbf{x}_n - \mathbf{m}_k\right)^\top.$$

• Чтобы обобщить внешнюю (межклассовую) дисперсию, просто возьмём остаток полной дисперсии

$$\begin{split} \mathbf{S}_T &= \sum_n \left(\mathbf{x}_n - \mathbf{m}\right) \left(\mathbf{x}_n - \mathbf{m}\right)^\top, \\ \mathbf{S}_B &= \mathbf{S}_T - \mathbf{S}_W. \end{split}$$

• Обобщить критерий можно разными способами, например:

$$J(\mathbf{W}) = \operatorname{Tr}\left[\mathbf{s}_W^{-1}\mathbf{s}_B\right],$$

где ${f s}$ – ковариации в пространстве проекций на ${f y}$:

$$\begin{split} \mathbf{s}_W &= \sum_{k=1}^K \sum_{n \in C_k} \left(\mathbf{y}_n - \boldsymbol{\mu}_k\right) \left(\mathbf{y}_n - \boldsymbol{\mu}_k\right)^\top, \\ \mathbf{s}_B &= \sum_{k=1}^K N_k \left(\boldsymbol{\mu}_k - \boldsymbol{\mu}\right) \left(\boldsymbol{\mu}_k - \boldsymbol{\mu}\right)^\top, \end{split}$$

где
$$\mu_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{y}_n$$
.

Спасибо!

Спасибо за внимание!