

# Заметки по многошаговому прогнозированию

Зехов Матвей

22 октября 2019 г.

## Содержание

<b>1 Основные стратегии многошагового прогнозирования</b>	<b>1</b>
1.1 Модель . . . . .	1
1.2 Рекурсивная стратегия . . . . .	3
1.3 Прямая стратегия . . . . .	5
<b>2 Комбинации прямой и рекурсивной стратегий</b>	<b>7</b>
2.1 DirRec . . . . .	7
2.2 MSRV . . . . .	7
2.3 Стратегия RECTIFY . . . . .	7
2.4 Оценка суммы сдвига и дисперсии . . . . .	8
2.5 Регрессионные методы и отбор моделей . . . . .	8
2.6 Оценка смещения и дисперсии . . . . .	8
<b>3 Стратегии с множественными горизонтами</b>	<b>9</b>
3.1 Прямые стратегии . . . . .	9
3.1.1 Стратегия DIRJOINT . . . . .	9
3.1.2 Стратегия SJOINT . . . . .	10
3.1.3 Стратегия DIRJOINTL . . . . .	10
3.1.4 Стратегия RECJOINT . . . . .	10
3.1.5 Стратегия RECJOINTL . . . . .	10
3.2 Рекурсивные стратегии . . . . .	10
3.2.1 Стратегия RECJOINT . . . . .	10
<b>4 Эксперименты и симуляции</b>	<b>11</b>
4.1 Параметры симуляций и результаты. . . . .	11

## 1 Основные стратегии многошагового прогнозирования

### 1.1 Модель

Общая идея всех подходов: Больше параметров  $\rightarrow$  точнее модель, более гибкая, меньше сдвиг, больше дисперсия. И наоборот.

Если нам важно более направление, чем значение прогноза, то следует выбирать модель с меньшим смещением, так как большое смещение искривляет направление. А если у нас короткие ряды, склонные к переобучению, то - с маленькой дисперсией.

Нам необходимо спрогнозировать временной ряд из  $T$  наблюдений на  $N$  шагов вперёд. Будем предполагать, что данные пришли к нам из модели (возможно нелинейной) следующего вида:

$$y_t = f(x_{t-1}) + \varepsilon_t \text{ with } x_t = [y_t, \dots, y_{t-d+1}]'$$

$\varepsilon_t$  - iid, среднее ноль, дисперсия  $\sigma^2$ ,  $\kappa = \mathbb{E}(\varepsilon^4) > 0$ .

Процесс специфицируется функцией  $f$ , размерностью эмбедингов  $d$  и ошибкой  $\varepsilon_t$ . Цель – оценить условное среднее  $\mu_{t+h|t} = \mathbb{E}(y_{t+h}|x_t)$ , и мы будем пытаться использовать различные стратегии аппроксимации  $\mu_{t+h|t}$ .

При горизонте прогнозирования один имеем просто  $\mu_{t+1|t} = f(x_t)$ . Если  $f$  линейная, можно записать более общую формулу.

$$\mu_{t+h|t} = \begin{cases} f\left([f^{(h-1)}(x_t), \dots, f^{(h-d)}(x_t)]'\right), & \text{if } h > 0 \\ x_t' w_h, & \text{if } 1-d \leq h \leq 0 \end{cases}$$

где  $w_h$  имеет единицу на  $j = 1 - h$  позиции и нули в иных. Если функция  $f$  линейна, то условное среднее можно посчитать рекурсивно. В иных случаях при  $h > 1$  и нелинейной  $f$  нет простой формы подсчёта  $\mu_{t+h|t}$ .

Каждая стратегия включает оценку одной или более моделей, которые не обязательно той же формы, что и  $f$  и могут иметь иную размерность эмбедингов. Для одношаговых прогнозов мы будем оценивать модель  $y_t = m(x_{t-1}; \theta) + e_t$  где  $x_t = [y_t, \dots, y_{t-p+1}]'$ . Мы будем оценивать форму  $m$ , параметры  $\theta$  и размерность  $p$ . Если форма  $m$  совпадает с формой  $f$ , то будем писать  $m \asymp f$ . В идеале мы бы хотели, чтобы  $p = d$ ,  $m \asymp f$  и  $\theta$  были близки к истинным параметрам, но мы не будем придерживаться этих предположений из-за возможных ошибок спецификации.

Обозначим  $\hat{m}^{(h)}(x_t)$  прогнозы конкретной стратегии на горизонте  $h$  и обозначим  $m^{(h)} = \mathbb{E}[\hat{m}^{(h)}(x_t) | x_t]$

Тогда MSE прогноза на горизонте  $h$  можно записать следующей формулой:

$$\text{MSE}_h = \mathbb{E} \left[ (y_{t+h} - \hat{m}^{(h)}(x_t))^2 \right] = \underbrace{\mathbb{E} \left[ (y_{t+h} - \mu_{t+h|t})^2 \right]}_{\text{Noise}} + \underbrace{(\mu_{t+h|t} - m^{(h)}(x_t))^2}_{\text{Bias } b_h^2} + \underbrace{\mathbb{E} \left[ (\hat{m}^{(h)}(x_t) - m^{(h)}(x_t))^2 \right]}_{\text{Variance}}$$

Декомпозиция получена следующим образом:

$$\begin{aligned} & \mathbb{E} \left[ (y - \hat{f})^2 \right] \\ &= \mathbb{E} \left[ (f + \varepsilon - \hat{f})^2 \right] \\ &= \mathbb{E} \left[ (f + \varepsilon - \hat{f} + \mathbb{E}[\hat{f}] - \mathbb{E}[\hat{f}])^2 \right] \\ &= \mathbb{E} \left[ (f - \mathbb{E}[\hat{f}])^2 \right] + \mathbb{E}[\varepsilon^2] + \mathbb{E} \left[ (\mathbb{E}[\hat{f}] - \hat{f})^2 \right] + 2\mathbb{E}[(f - \mathbb{E}[\hat{f}])\varepsilon] + 2\mathbb{E}[\varepsilon(\mathbb{E}[\hat{f}] - \hat{f})] \\ &\quad + 2\mathbb{E}[(\mathbb{E}[\hat{f}] - \hat{f})(f - \mathbb{E}[\hat{f}])] \\ &= (f - \mathbb{E}[\hat{f}])^2 + \mathbb{E}[\varepsilon^2] + \mathbb{E} \left[ (\mathbb{E}[\hat{f}] - \hat{f})^2 \right] + 2(f - \mathbb{E}[\hat{f}])\mathbb{E}[\varepsilon] + 2\mathbb{E}[\hat{f}] - \hat{f} + 2\mathbb{E}[\mathbb{E}[\hat{f}] - \hat{f}](f - \mathbb{E}[\hat{f}]) \\ &= (f - \mathbb{E}[\hat{f}])^2 + \mathbb{E}[\varepsilon^2] + \mathbb{E} \left[ (\mathbb{E}[\hat{f}] - \hat{f})^2 \right] \\ &= (f - \mathbb{E}[\hat{f}])^2 + \text{Var}[y] + \text{Var}[\hat{f}] \\ &= \text{Bias}[\hat{f}]^2 + \text{Var}[y] + \text{Var}[\hat{f}] \\ &= \text{Bias}[\hat{f}]^2 + \sigma^2 + \text{Var}[\hat{f}] \end{aligned}$$

Более того, сдвиг можно разложить ещё на две составляющие:

$$\begin{aligned}
B_h &= \mathbb{E}_{x_t} \left[ \left( \mu_{t+h|t} - \mathbb{E}_{Y_T} \left[ m \left( x_t; \hat{\theta}_{Y_T}; h \right) \right] \right)^2 \right] \\
&= \mathbb{E}_{x_t} \left[ \underbrace{\left( \mu_{t+h|t} - m(x_t; \theta^*; h) \right)}_A + \underbrace{m(x_t; \theta^*; h) - \mathbb{E}_{Y_T} \left[ m \left( x_t; \hat{\theta}_{Y_T}; h \right) \right]}_B \right]^2
\end{aligned}$$

В этой части сознательно использованы несколько другие, более точные обозначения, которые позволят разобраться полнее. Интуитивно понять аналогию несложно. Часть А отображает различие между условным средним семейства моделей, которое мы рассматриваем. Например, если первое - нелинейное, а мы оцениваем линейную модель. Тогда правая часть А будет наилучшей из достижимых оценкой линейных параметров. Следует обратить внимание, что подставлены именно оптимальные параметры. Тогда А обозначает ограничения выбранной нами модели относительно истинной зависимости.

Часть В отображает ошибку вследствие ограниченности временных рядов до Т наблюдений. Она выражает как конечность выборки влияет на сдвиг. Это можно понять по тому, что в левой части В подставлены оптимальные параметры, а в правой - оценённые по конечному ряду.

Следовательно, даже если мы подберём оптимальную модель, то мы не сможем избавиться от части В вследствие оконечной выборки.

Важно заметить, что компонента шума не зависит ни от рассматриваемой модели, ни от стратегии построения прогноза. Она зависит только от процесса генерации данных.

При росте числа наблюдений дисперсия будет сходиться к нулю, и далее она исключается из рассмотрения. Далее разложим левую часть (Noise) для горизонта  $h = 2$ . Аналогично для любого другого.

$$y_{t+2} = f(y_{t+1}, \dots, y_{t-d+2}) + \varepsilon_{t+2} = f(f(x_t) + \varepsilon_{t+1}, \dots, y_{t-d+2}) + \varepsilon_{t+2}$$

Используя разложение Тейлора до второго порядка (как я понял, относительно переменной  $\varepsilon_t$ ), получим:

$$y_{t+2} \approx f(f(x_t), \dots, y_{t-d+2}) + \varepsilon_{t+1} f_{x_1} + \frac{1}{2} (\varepsilon_{t+1})^2 f_{x_1 x_1} + \varepsilon_{t+2},$$

где  $f_{x_1}$  - первая производная  $f$  по её первому аргументу, а  $f_{x_1 x_1}$  - вторая производная дважды по первому аргументу

Тогда шум можно раскрыть по следующей формуле (правая часть это математическое ожидание левой во второй строке):

$$\begin{aligned}
&\mathbb{E} \left[ (y_{t+2} - \mu_{t+2|t})^2 \right] \\
&\approx \mathbb{E} \left[ \left( f(f(x_t), \dots, y_{t-d+2}) + \varepsilon_{t+1} f_{x_1} + \frac{1}{2} \varepsilon_{t+1}^2 f_{x_1 x_1} + \varepsilon_{t+2} - f(f(x_t), \dots, y_{t-d+2}) - \frac{1}{2} \sigma^2 f_{x_1 x_1} \right)^2 \right] \\
&= \sigma^2 (1 + f_{x_1}^2) + \frac{1}{4} (\kappa - \sigma^4) f_{x_1 x_1}^2
\end{aligned}$$

Тогда MSE на горизонте  $h = 2$  будет равно:

$$\text{MSE}_2 \approx \sigma^2 (1 + f_{x_1}^2) + \frac{1}{4} (\kappa - \sigma^4) f_{x_1 x_1}^2 + (\mu_{t+h|t} - m^{(h)}(x_t))^2$$

## 1.2 Рекурсивная стратегия

При рекурсивной стратегии оценивается модель

$$y_t = m(x_{t-1}; \theta) + e_t, \quad \text{where } x_t = [y_t, \dots, y_{t-p+1}]'$$

где  $\mathbb{E}[e_t] = 0$ , целевая функция:  $\mathbb{E}[(y_{t+1} - m(\mathbf{x}_t; \boldsymbol{\theta}))^2 | \mathbf{x}_t]$ , а параметры оцениваются следующим образом:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} \sum_t (y_t - m(\mathbf{x}_{t-1}; \boldsymbol{\theta}))^2$$

Прогнозы вычисляются рекурсивно относительно предыдущих:

$$\hat{m}^{(h)}(x_t) = \begin{cases} \hat{m} \left( [\hat{m}^{(h-1)}(x_t), \dots, \hat{m}^{(h-p)}(x_t)]' \right), & \text{if } h > 0 \\ x_t' w_h, & \text{if } 1 - p \leq h \leq 0 \end{cases}$$

где  $\hat{m}(x) = m(x; \hat{\boldsymbol{\theta}})$ . Иногда это называют итеративным многошаговым прогнозированием.

Элемент смещения  $\delta(z_t; \boldsymbol{\theta})$  возникает из-за:

- ☀ Недостатка гибкости рассматриваемой модели.
- ☀ Потенциального пропуска регрессоров во входных данных
- ☀ Неадекватного алгоритма оценки параметров  $\theta$ .

Вариативность прогноза описывается элементом  $\eta(z_t; \boldsymbol{\theta}) \varepsilon_\eta$ . Она возникает вследствие следующих причин:

- ☀ Ограниченности длины ряда при оценке  $\theta$
- ☀ Включение в  $z$  потенциально избыточных или бессмысленных переменных.
- ☀ Избыточная гибкость модели и переобучение

Предсказания рекурсивной стратегии:

$$m_h(x_t; \hat{\boldsymbol{\theta}}) = \underbrace{\underbrace{f(\mathbf{x}_t)}_{\mu_{t+h|t}} + \delta(z_t; \boldsymbol{\theta}) + \eta(z_t; \boldsymbol{\theta}) \varepsilon_\eta}_{\underbrace{m(z_t; \boldsymbol{\theta})}_{m(z_t; \hat{\boldsymbol{\theta}})}}$$

Выбор квадратичного функционала для  $\hat{\boldsymbol{\theta}}$  обеспечивает  $m^{(1)}(x_t) = \mu_{t+1|t}$ , и, как следствие,  $m \asymp f$  и  $p = d$ . Следовательно, одношаговый прогноз будет несмещённым. Однако это не сохраняется для порядка 2 и более:

$$\begin{aligned} b_2 = \mu_{t+2|t} - m^{(2)}(\mathbf{x}_t) &\approx f(f(\mathbf{x}_t), \dots, y_{t-d+2}) + \frac{1}{2} \sigma^2 f_{x_1 x_1} - m(m(\mathbf{x}_t), \dots, y_{t-p+2}) \\ &\approx [f(f(\mathbf{x}_t), \dots, y_{t-d+2}) - m(m(\mathbf{x}_t), \dots, y_{t-p+2})] + \frac{1}{2} \sigma^2 f_{x_1 x_1} \end{aligned}$$

Можно аналогично раскрыть и форму:

$$\begin{aligned} m(x_t; \hat{\boldsymbol{\theta}}; 2) &= m(m(x_t; \hat{\boldsymbol{\theta}}), \dots, y_{t-p+2}; \hat{\boldsymbol{\theta}}) \approx f(f(x_t), \dots, y_{t-p+2}) \\ &+ \delta(f(x_t), \dots, y_{t-p+2}; \boldsymbol{\theta}) + \eta(f(x_t), \dots, y_{t-p+2}; \boldsymbol{\theta}) \varepsilon_{\eta_2} \\ &+ \delta(z_t; \boldsymbol{\theta}) m_{z_1} + \frac{1}{2} [\delta(z_t; \boldsymbol{\theta})]^2 m_{z_1 z_1} + \eta(z_t; \boldsymbol{\theta}) \varepsilon_{\eta_1} m_{z_1} + \frac{1}{2} [\eta(z_t; \boldsymbol{\theta}) \varepsilon_{\eta_1}]^2 m_{z_1 z_1} \end{aligned}$$

Следовательно, даже при идеальных условиях  $m \asymp f$  и  $p = d$ , смещение будет отсутствовать только при линейности истинной функции  $f$ .

Для этой стратегии MSE будет равна:

$$\text{MSE}_2^{\text{recursive}} \approx \sigma^2 (1 + f_{x_1}^2) + \frac{1}{4} (\kappa - \sigma^4) f_{x_1 x_1}^2 + ([f(f(x_t), \dots, y_{t-d+2}) - m(m(x_t), \dots, y_{t-p+2})] + \frac{1}{2} \sigma^2 f_{x_1 x_1})^2$$

При выполнении условий  $m \asymp f$  и  $p = d$  формула упрощается до

$$\text{MSE}_2^{\text{recursive}} \approx \sigma^2 (1 + f_{x_1}^2) + \frac{1}{4} \kappa f_{x_1 x_1}^2$$

Когда модель неправильно специфицирована или параметры оценены неверно, ашибка существенно больше.

Есть различные вариации:

С оценкой своего параметра для каждого горизонта:

$$\hat{\theta}_h = \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_t [y_t - m^{(h)}(x_{t-h}; \theta)]^2$$

Такую стратегию назовём RECMULTI

Минимум по первым  $H$  горизонтам:

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_{h=1}^H \sum_t [y_t - m^{(h)}(x_{t-h}; \theta)]^2$$

Достоинства: мало вычислений, только одна модель, низкая дисперсия, может быть лучше прямого метода, если модель правильно специфицирована

Недостатки: Смещение, только одна модель

В случае AR(p) гауссовского-процесса генерации данных и AR(k) оцениваемой модели, то при  $k \geq p$  можно показать, что оценки МНК асимптотически эквивалентны ML и все оценки с индексом  $i > p = 0$ . В общем, параметры прогноза, которые зависят от параметров оценённой модели по инвариантности оценок ML можно тоже оценить.

$$\hat{\mathbf{a}}_n(h, k) = \hat{A}_n^{h-1}(k) \hat{\mathbf{a}}_n(1, k)$$

is asymptotically equivalent to the MLE of

$$\mathbf{a}(h, k) = (a_1(h, k), \dots, a_k(h, k))' = A^{h-1}(k) \mathbf{a}(1, k)$$

where  $k \geq p$ ,  $A^0(k) = I_k$ , and

$$A(k) = \left( \mathbf{a}(1, k) \middle| \frac{I_{k-1}}{\mathbf{0}_{k-1}'} \right)$$

### 1.3 Прямая стратегия

Прямая стратегия: давайте оценивать для каждого горизонта свою модель.

$$y_t = m_h(y_{t-h}, \dots, y_{t-h-p_h}; \theta_h) + e_{t,h}$$

Тогда параметры для каждой модели будем оценивать МНК по соответствующим пространствам параметров и  $\hat{m}^{(h)}(x_t) = m_h(x_t; \hat{\theta}_h)$

$$\hat{\theta}_h = \underset{\theta_h \in \Theta_h}{\operatorname{argmin}} \sum_t [y_t - m_h(x_{t-h}; \theta_h)]^2$$

Такой подход обычно называют "direct multi-step". Для оценки  $h$  моделей требуется больше времени. Теперь мы оцениваем не единую модель, а несколько независимых моделей

Из-за квадратичных отклонений и  $m^{(h)}(x_t) = m_h(x_t)$  (истинная модель для горизонта  $h$  это  $h$ -ая модель) bias модели при горизонте 2 равен:

$$b_2 = \mu_{t+2|t} - m^{(2)}(x_t) \approx f(f(x_t), \dots, y_{t-d+2}) + \frac{1}{2}\sigma^2 f_{x_1 x_1} - m_2(y_t, \dots, y_{t-p+1})$$

Тогда эта стратегия буде вести к несмещённым оценкам когда  $m_2(y_t, \dots, y_{t-p+1}) \asymp f(f(x_t), \dots, y_{t-d+2}) + \frac{1}{2}\sigma^2 f_{x_1 x_1}$ . Эти условия будут выполняться в случае, если  $m_2$  будет достаточно гибкой моделью.

Можно также дополнительно раскрыть формулу прогноза следующим образом:

$$m(x_t; \hat{\theta}_2; 2) = \underbrace{\mu_{t+2|t} + \delta(r_t; \theta_2)}_{m_2(r_t; \theta_2)} + \eta(r_t; \theta_2) \varepsilon_\eta$$

Следует обратить внимание, что в отличие от рекурсивной стратегии, в прямой сразу может возникнуть условное среднее без рекурсивного вызова функции от  $x_t$

Во всех иных стратегиях прогноз будет выглядеть следующим образом:

$$g(x_t; \hat{\gamma}; 2) = \underbrace{\mu_{t+2|t} + \delta(r_t; \gamma)}_{m_2(r_t; \gamma)} + \eta(r_t; \gamma) \varepsilon_\eta$$

где гамма, это набор параметров, различный для каждой стратегии. Все стратегии за исключение рекурсивной используют  $h$ -шаговую ошибку в качестве целевой функции.

Вспоминая формулу

$$\text{MSE}_2 \approx \sigma^2 (1 + f_{x_1}^2) + \frac{1}{4} (\kappa - \sigma^4) f_{x_1 x_1}^2 + (\mu_{t+h|t} - m^{(h)}(x_t))^2$$

можем получить аналогичную для прямого случая

$$\text{MSE}_2^{\text{direct}} \approx \sigma^2 (1 + f_{x_1}^2) + \frac{1}{4} (\kappa - \sigma^4) f_{x_1 x_1}^2 + \left[ f(f(x_t), \dots, y_{t-d+2}) + \frac{1}{2}\sigma^2 f_{x_1 x_1} - m_2(y_t, \dots, y_{t-p+1}) \right]^2$$

Если истинная стратегия несмещённая, то уравнение упростится до

$$\text{MSE}_2^{\text{direct}} \approx \sigma^2 (1 + f_{x_1}^2) + \frac{1}{4} (\kappa - \sigma^4) f_{x_1 x_1}^2$$

В отличие от рекурсивного случае здесь  $m_2$  отражает не всю истинную модель и может быть достаточно гибкой чтобы подстроиться. Следовательно, несмещённый случай возможен. не только в линейном случае.

Получается, в идеальных условиях, когда  $m \asymp f$  and  $p = d$  для рекурсивной стратегии и прямая стратегия несмещённая, можно посчитать разницу ошибок:

$$\text{MSE}_2^{\text{recursive}} - \text{MSE}_2^{\text{direct}} \approx \frac{1}{4} \sigma^4 f_{x_1 x_1}^2$$

Были работы, которые показали, что даже при линейности всех истинных зависимостей и стационарности ряда рекурсивная модель хуже, чем прямая.

У прямого метода помимо вычислительной сложности есть ещё проблема построения слишком разных моделей на разных горизонтах. Следовательно прогнозы строятся на независимых моделях с различными формами и условной информацией. Это как раз и увеличивает вариацию прогноза, но зато смещение низкое на нелинейных данных.

Может работать лучше рекурсивной стратегии, если она модель неверно специфицирована, так как ошибки спецификации обычно создают сдвиг. Ещё один недостаток: автокорреляция ошибок.

## 2 Комбинации прямой и рекурсивной стратегий

### 2.1 DirRec

Эта смешанная стратегия использует различные наборы параметров  $\theta_h$  для каждого горизонта  $h$ , как в прямой стратегии, но включает предыдущие прогнозы вместе с другими входными значениями:

$$\hat{\theta}_h = \operatorname{argmin}_{\theta_h \in \Theta_h} \sum_t [y_t - [m_h(\hat{m}_{h-1}, \dots, \hat{m}_1, \mathbf{r}_{t-h}; \theta_h)]]^2$$

где  $\hat{m}_h$  – это сокращение для  $m_h(\mathbf{r}_{t-h}; \hat{\theta}_h)$

Тогда прогнозы для каждого горизонта из соответствующей модели получаются следующим образом:  $\hat{\mu}_{T+h|T} = m_h(\hat{m}_{h-1}, \dots, \hat{m}_1, \mathbf{r}_T; \hat{\theta}_h)$

### 2.2 MSRV

Ещё одна комбинация прямого и рекурсивного подходов. Если мы представим горизонты прогнозирования следующим образом:  $H = L \times R$ , то эта стратегия первые  $L$  моделей оценит как прямые, а потом использует их  $R$  раз для получения  $H$  прогнозов:

$$\hat{\mu}_{T+h|T} = \begin{cases} m_l(\mathbf{r}_T; \hat{\theta}_l) & \text{if } h \leq L \\ m_l(\hat{m}_{l-1}, \dots, \hat{m}_1, \mathbf{r}'_T; \hat{\theta}_l) & \text{if } h > L \end{cases}$$

где  $l = (h - 1) \% L + 1$ ,  $\hat{m}_l$  это сокращение для  $m_l(\mathbf{r}_T; \hat{\theta}_l)$ , и  $\mathbf{r}'_T \subset \mathbf{r}_T$

Преимущество: время вычислений. Однако, из-за рекурсии точность также страдает.

### 2.3 Стратегия RECTIFY

Основная идея - взять плюсы двух первых. Она начинается с рекурсивных прогнозов и улучшает их так, что они становятся несмещёнными и имеют меньшую ошибку.

Начнём с рекурсивных прогнозов на линейной модели. Они будут смещёнными даже когда процесс правильно специфицирован. Потом скорректируем их моделированием прогнозных ошибок через прямую стратегию и получим несмещённые прогнозы.

Достоинство подхода в том, что он соединяет все модели прямого подхода через единую базовую модель, снижая разрегулировку, порождённую независимыми моделями. Грубо говоря, прямые модели станут более похожими друг на друга.

Обозначим базовую линейную модель как  $y_t = z(x_{t-1}; \theta) + e_t$ . Из неё получим рекурсивные прогнозы. После этого мы изменим прогнозы базовой модели, применяя модели прямого прогнозирования к ошибкам рекурсивного прогноза. То есть мы будем фитить следующие модели:

$$y_t - z^{(h)}(y_{t-h}, \dots, y_{t-h-p}; \hat{\theta}) = r_h(y_{t-h}, \dots, y_{t-h-p_h}; \gamma_h) + e_{t,h}, \forall h = 1 \dots H$$

Все параметры оценены через OLS. Прогнозы получены через суммирование базовой и дополнительной модели:  $\hat{m}^{(h)}(\mathbf{x}_t) = \hat{z}^{(h)}(\mathbf{x}_t) + \hat{r}_h(\mathbf{x}_t)$

Пусть  $m^{(h)}(\mathbf{x}_t) = \mathbb{E}[\hat{m}^{(h)}(\mathbf{x}_t) | \mathbf{x}_t]$ . Тогда смещение можно записать как

$$\begin{aligned}
b_2 &= \mu_{t+2|t} - m^{(2)}(\mathbf{x}_t) \\
&\approx f(f(\mathbf{x}_t), \dots, y_{t-d+2}) + \frac{1}{2}\sigma^2 f_{x_1 x_1} - [z(z(\mathbf{x}_t), \dots, y_{t-p+2}) + r_2(y_t, \dots, y_{t-p_2+1})] \\
&= \left[ f(f(\mathbf{x}_t), \dots, y_{t-d+2}) - z(z(\mathbf{x}_t), \dots, y_{t-p+2}) + \frac{1}{2}\sigma^2 f_{x_1 x_1} \right] - r_2(y_t, \dots, y_{t-p_2+1})
\end{aligned}$$

Следовательно, стратегия будет несмещённой, когда  $r_2(y_t, \dots, y_{t-p_2+1}) \asymp f(f(x_t), \dots, y_{t-d+2}) - z(z(x_t), \dots, y_{t-p+2}) + \frac{1}{2}\sigma^2 f_{x_1 x_1}$ , то есть когда дополнительная модель достаточно гибкая, чтобы оценить условное среднее ошибок базовой модели.

Смещение базовой модели скорректировано дополнительной моделью. Следовательно, отпадает необходимость условия  $z \asymp f$ . В случае статьи рассматривалась модель линейной авторегрессии с отбором лага по AIC. Конечно, модель будет смещённой для нелинейных моделей, но она поможет моделировать большую часть сигнала  $f$  и будет давать относительно маленькую дисперсию из-за своей линейности. Дополнительная модель должна быть достаточно гибкой и в статье рассматриваются оценки KNN для неё.

Можно как и ранее вывести MSE для этой стратегии.

$$\begin{aligned}
\text{MSE}_2^{\text{rectify}} &\approx \sigma^2 (1 + f_{x_1}^2) + \frac{1}{4} (\kappa - \sigma^4) f_{x_1 x_1}^2 \\
&\quad + \left[ \left[ f(f(x_t), \dots, y_{t-d+2}) - z(z(x_t), \dots, y_{t-p+2}) + \frac{1}{2}\sigma^2 f_{x_1 x_1} \right] - r_2(y_t, \dots, y_{t-p_2+1}) \right]^2
\end{aligned}$$

Когда стратегия несмещённая, она имеет такую же MSE, как и прямая стратегия.

$$\text{MSE}_2^{\text{recursive}} - \text{MSE}_2^{\text{rectify}} = \text{MSE}_2^{\text{recursive}} - \text{MSE}_2^{\text{direct}} \approx \frac{1}{4}\sigma^4 f_{x_1 x_1}^2$$

Однако, надо не забывать учитывать дисперсию. Хотя у обеих стратегий они стремятся к нулю, у модели rectify меньшая дисперсия на конечных выборках.

## 2.4 Оценка суммы сдвига и дисперсии

Для горизонта  $h = 2$  можно записать следующую формулу:

$$B_2(\mathbf{x}_t) + V_2(\mathbf{x}_t) \approx (\mu_{t+2|t} - m(z_t; \boldsymbol{\theta}; 2))^2 + \mathbb{E}_{Y_T} \left[ \left( m(z_t; \hat{\boldsymbol{\theta}}; 2) - m(z_t; \boldsymbol{\theta}; 2) \right)^2 | \mathbf{x}_t \right]$$

## 2.5 Регрессионные методы и отбор моделей

В статье - линейная авторегрессия и KNN.

Линейная модель отбирает порядок  $p$  по AIC.

KNN - нелинейная и непараметрическая модель. Использовалось взвешенное среднее.  $k$  - количество взвешенно усреднённых соседей. Чем больше соседей - тем больше сглаживание, меньше дисперсия, больше смещение и наоборот

## 2.6 Оценка смещения и дисперсии

$$\begin{aligned}
\text{MSE}_h &= \frac{1}{LR} \sum_{i=1}^L \sum_{j=1}^R \left( y_j(h) - \hat{m}_{D^i}^{(h)}(\mathbf{x}_j) \right)^2 \\
&= \text{Noise}_h + \text{Bias}_h^2 + \text{Variance}_h
\end{aligned}$$



Модели генерации данных:

☀ AR(6) - процесс,  $\varepsilon_t \sim \text{NID}(0,1)$

$$y_t = 1.32y_{t-1} - 0.52y_{t-2} - 0.16y_{t-3} + 0.18y_{t-4} - 0.26y_{t-5} + 0.19y_{t-6} + \varepsilon_t$$

☀ STAR,  $\varepsilon_t \sim \text{NID}(0, \sigma^2)$ ,  $\sigma^2 = [0.05^2, 0.1^2]$

$$y_t = 0.3y_{t-1} + 0.6y_{t-2} + (0.1 - 0.9y_{t-1} + 0.8y_{t-2}) [1 + e^{(-10y_{t-1})}]^{-1} + \varepsilon_t$$

### 3 Стратегии с множественными горизонтами

Все стратегии, описанные выше, можно отнести к стратегиям с одним горизонтом, то есть стратегиям, когда модель рассматривает каждый горизонт отдельно. Другой тип моделей будет рассматривать несколько горизонтов за один приём. Целевая функция для оценки параметров будет учитывать одновременно ошибки нескольких горизонтов и, соответственно, оценивать один набор параметров.

Идея - учесть общие характеристики моделей нескольких различных горизонтов вследствие автокорреляции.

Учёт параметров в такой присоединённой манере позволит:

- ☀ Учесть взаимозависимости между различными горизонтами прогнозирования и улучшить обобщающую способность модели.
- ☀ Избежать разрегулировки вследствие использования слишком разных моделей на каждом горизонте
- ☀ Это помогает при маленьком размере семпла, используя дополнительные семплы из связанных задач ??

Для получения прогноза мультигоризонтной модели, сравнимого с одногоризонтной, нужно чтобы на горизонте  $h$

#### 3.1 Прямые стратегии

Dirjoint и Recjoint дают больший вес поздним горизонтам, так как усредняют ошибки. Ошибки дальних горизонтах обычно сильно выше. Следовательно, возникает ещё смещение на ближних горизонтах. Dirjoint всё же лучше по смещению.

##### 3.1.1 Стратегия DIRJOINT

Эта стратегия передаёт параметры между всеми горизонтами и оценивает параметры по средней ошибке на всех горизонтах.

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \sum_t \frac{1}{H} \sum_{h=1}^H [y_t - m(r_{t-h}; \theta)]^2$$

Эта стратегия применялась в контексте нейронных сетей. Очень легко понять на примере простого многослойного пресептрона, в котором входной слой имеет размерность  $p$ , а выходной - размерность  $H$ . Прогнав все наблюдения на нём и минимизировав ошибку как раз и получим формулу выше. Если есть аналитическое решение в разумный срок, то ок, а иначе - прошагаем градиентным спуском.

### 3.1.2 Стратегия SJOINT

Эта модель также предполагает разделение моделей на  $L$  различных многогоризонтных моделей, в которых параметры оцениваются следующим образом:

$$\hat{\theta}_l = \operatorname{argmin}_{\theta_l \in \Theta_l} \sum_t \frac{1}{R} \sum_{h=(l-1)R+1}^{l \times R} [y_t - m_l(\mathbf{r}_{t-h}; \theta_l)]^2$$

Горизонты группируются через параметр  $l$ , и для каждого из них будет свой набор параметров  $\theta_l$

### 3.1.3 Стратегия DIRJOINTL

Эта модель усредняет ошибки на текущем горизонте,  $i$  предыдущих горизонтах и  $j$  последующих горизонтах. Аналогично, можно использовать в нейросетях.

$$\hat{\theta}_h = \operatorname{argmin}_{\theta_h \in \Theta_h} \sum_t \frac{1}{i+j+1} \sum_{h'=h-i}^{h+j} [y_t - m_h(\mathbf{r}_{t-h'}; \theta)]^2$$

Как по ощущениям, все эти модели пытаются так или иначе сохранить несмещённость, но уменьшить дисперсию, уйдя от количества параметров. Например, DIRJOINT уже начинает сильно напоминать рекурсивную стратегию. Разве что усреднение по горизонтам может дать чуть большую устойчивость. Да и то слегка сомнительно.

### 3.1.4 Стратегия RECJOINT

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \sum_t \frac{1}{H} \sum_{h=1}^H [y_t - m^{(h)}(z_{t-h}; \theta)]^2$$

Различие с обычной рекурсивной стратегией в том, что здесь при оптимизации параметров учитывается влияние всех горизонтов. Эта штука очень похожа на BPTT из LSTM.

### 3.1.5 Стратегия RECJOINTL

Аналогично можем усреднять не все горизонты, а плюс-минус вперёд и назад.

$$\hat{\theta}_h = \operatorname{argmin}_{\theta_h \in \Theta_h} \sum_t \frac{1}{i+j+1} \sum_{h'=h-i}^{h+i} [y_t - m^{(h)}(\mathbf{r}_{t-h'}; \theta)]^2$$

Количество горизонтов, включённых в какую-либо оптимизируемую функцию далее будем обозначать как  $L_h \subseteq \{1 \dots, H\}$ . В случае, например, DIRJOINTL:  $L_h = \{h-i, \dots, h, \dots, h+i\}$

## 3.2 Рекурсивные стратегии

### 3.2.1 Стратегия RECJOINT

Логичная идея: а давайте также усредним по всем горизонтам, но уже рекурсивную стратегию в рамках модели, прогнозирующей рекурсивно.

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \sum_t \frac{1}{H} \sum_{h=1}^H [y_t - m^{(h)}(z_{t-h}; \theta)]^2$$

Аналогично, можно усреднять не все горизонты, а несколько вперёд и назад:

$$\hat{\theta}_h = \operatorname{argmin}_{\theta_h \in \Theta_h} \sum_t \frac{1}{i+j+1} \sum_{h'=h-i}^{h+i} [y_t - m^{(h)}(\mathbf{r}_{t-h'}; \theta)]^2$$

## 4 Эксперименты и симуляции

### 4.1 Параметры симуляций и результаты.

На длинных рядах прямая стратегия превосходит рекурсивную и имеет малый сдвиг.

Хорошая практика – обрезать первые три тысячи сгенерированных наблюдений для стабилизации временного ряда.

На малых горизонтах прогнозирования рекурсивная и прямая стратегия очень похожи.

Прямые прогнозы постепенно сходятся к средним прогнозам с ростом длины ряда и горизонта прогнозирования.

Нелинейность ряда и нелинейность модели прогнозирования ухудшает ухудшают сдвиг в рекурсивной стратегии. Рассматривая KNN-KNN можно заметить, что сильно уменьшается сдвиг, но сильно увеличивается дисперсия во всех моделях. Следовательно, не лучшая идея брать нелинейную модель в качестве базовой. В данном случае LINARP-KNN работала лучше. Особенно заметно на коротком горизонте. В сравнении с прямой стратегией этот вариант сильно лучше прямой стратегии на коротких горизонтах и примерно сопоставим на длинных. Также она сильно лучше на коротких рядах и приближается к прямой стратегии с ростом длины ряда.

В принципе, использование rectify всегда может быть оправданным чтобы избежать выбора между рекурсивной и прямой стратегиями.